

# Deep Learning

If not stated otherwise graphs and code examples were taken with permission from [Michael Nielsen](#)'s online book [Neural Networks and Deep Learning](#).

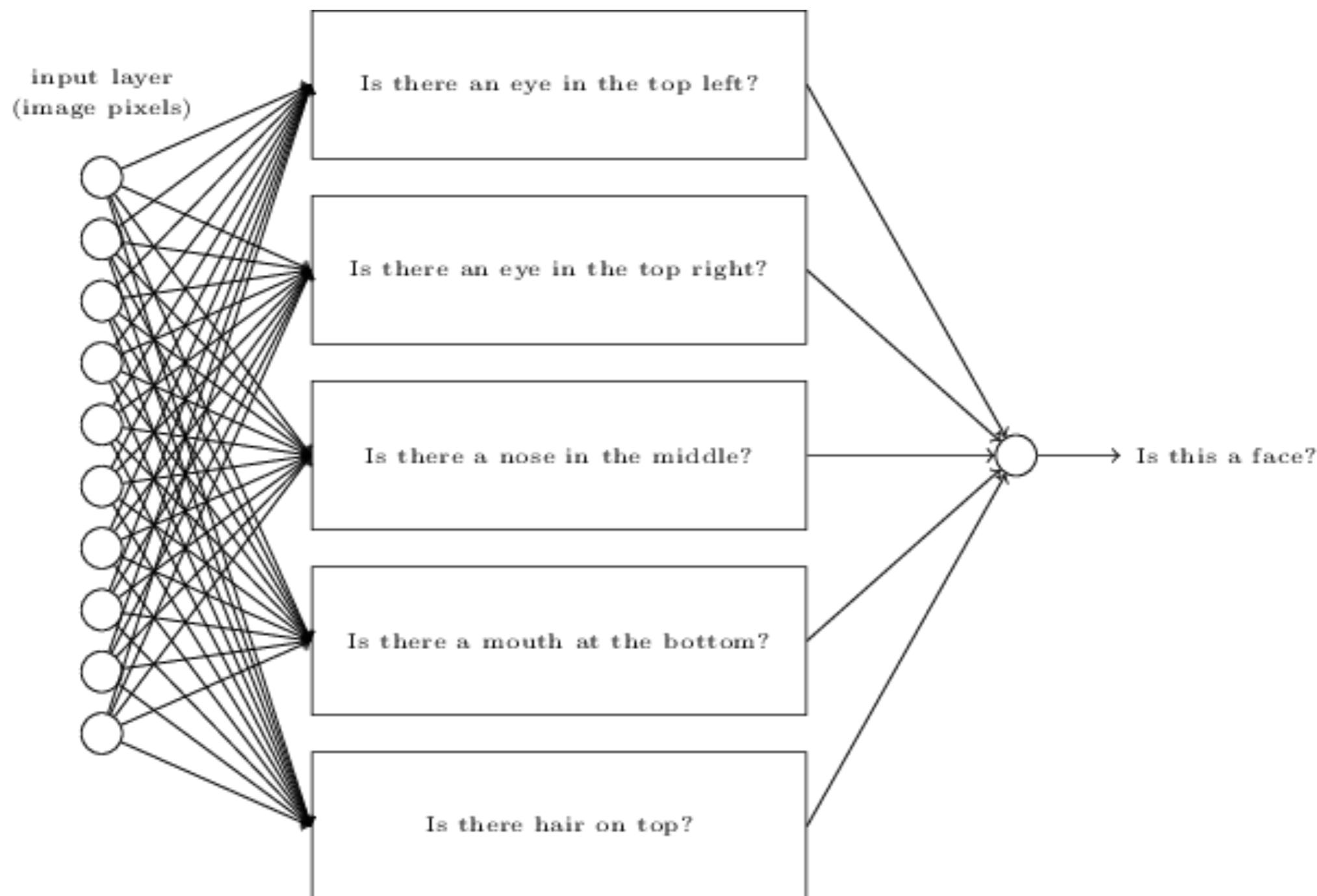
# Deep learning

- is a host of statistical machine learning techniques
- is generally based on artificial neural networks
- enables the automatic learning of feature hierarchies

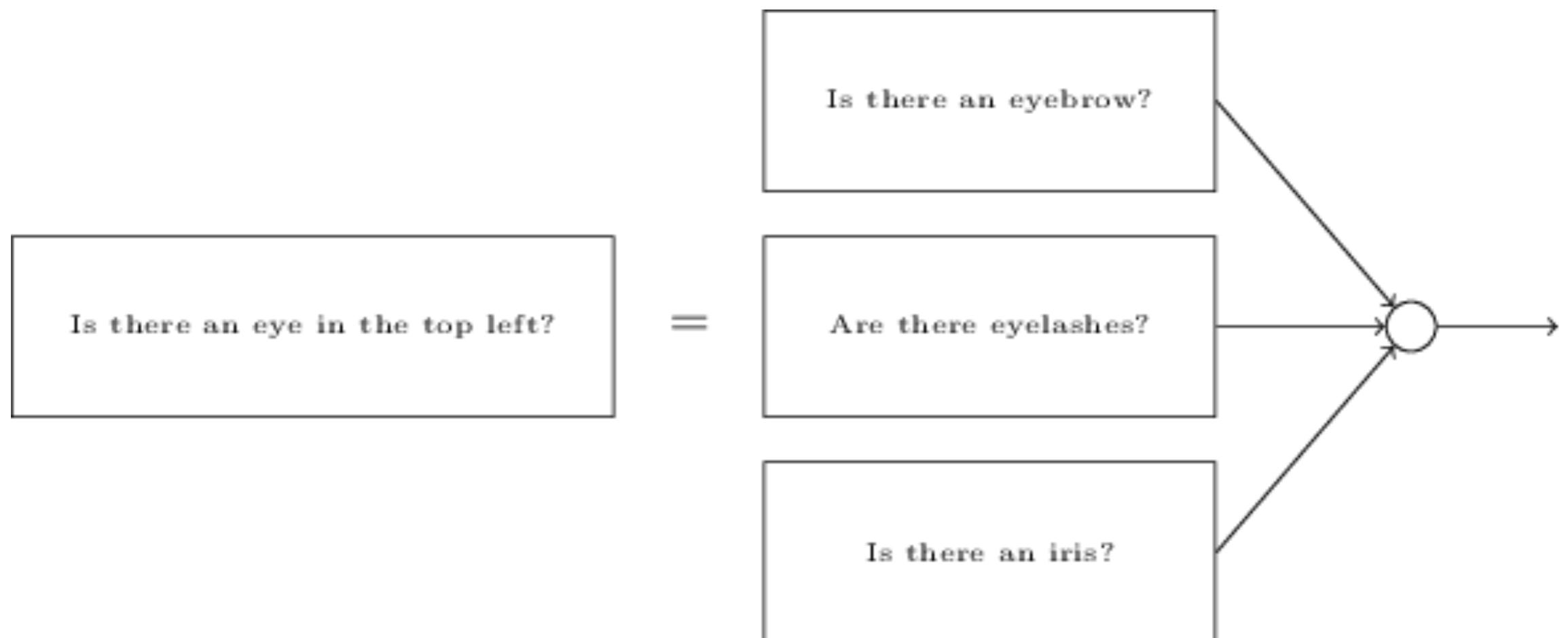
# Deep learning

- is a host of statistical machine learning techniques
- is generally based on artificial neural networks
- enables the automatic learning of **feature hierarchies**

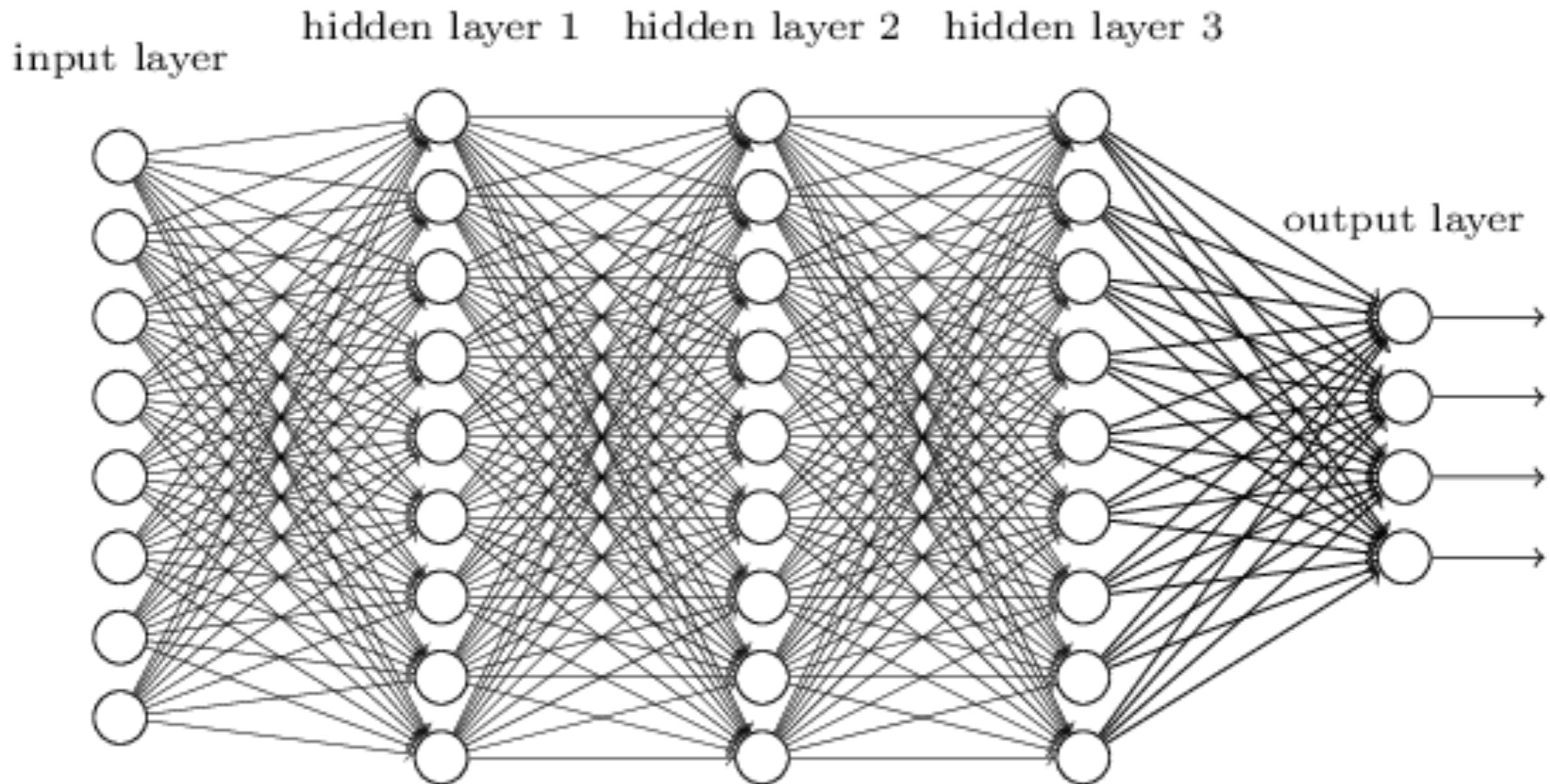
# Feature hierarchy



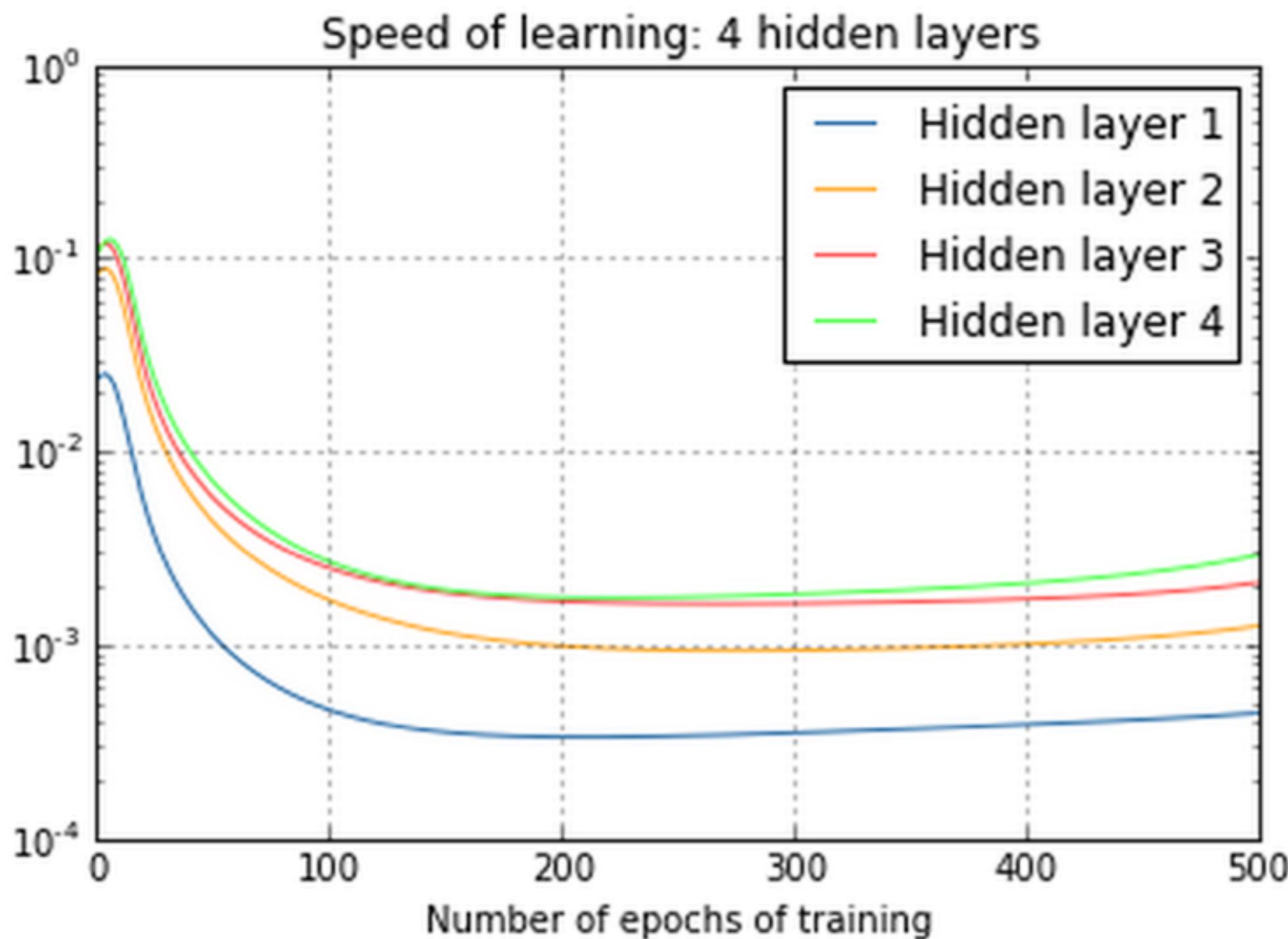
# Feature hierarchy



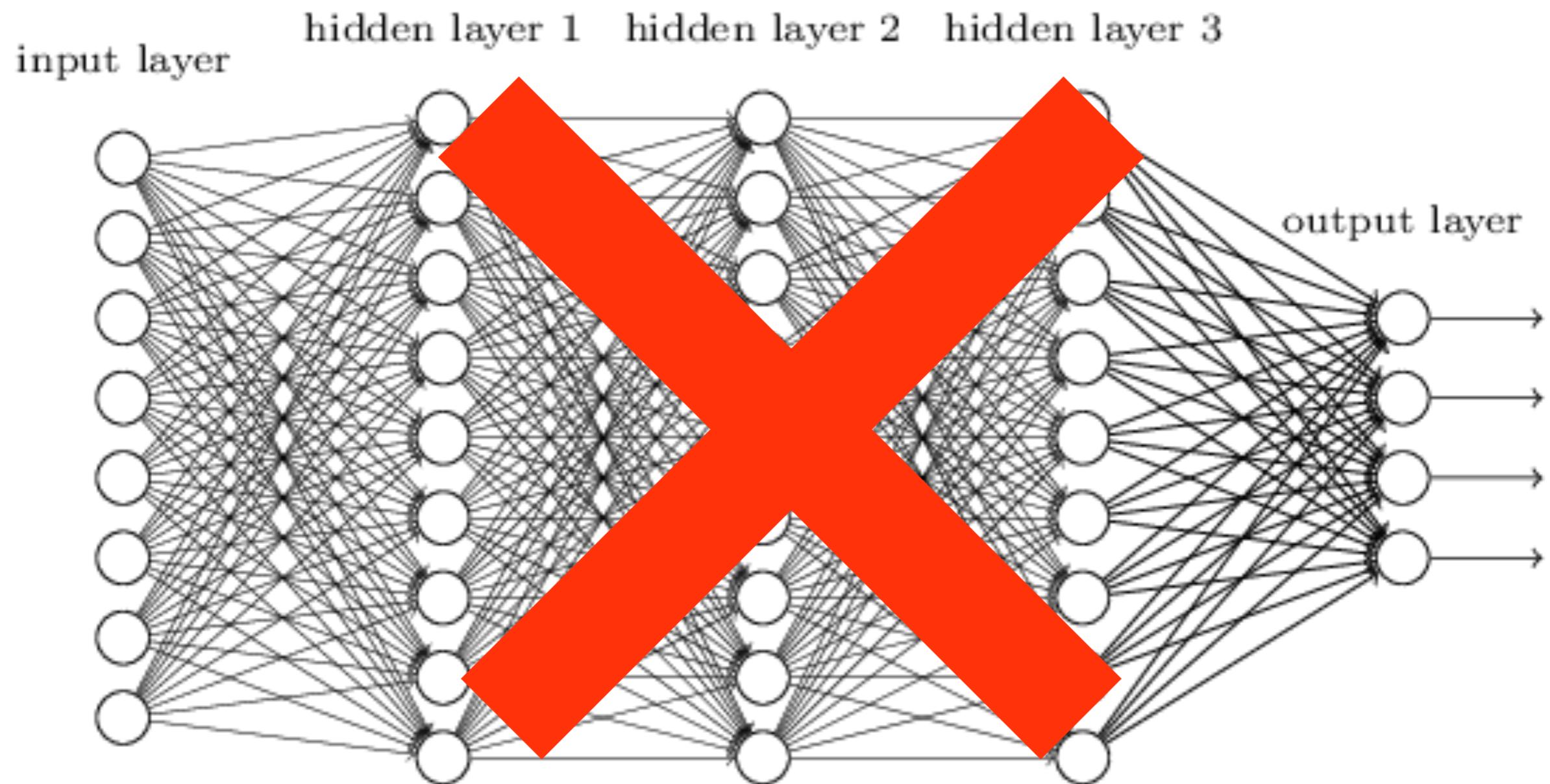
# Let's train our first deep network



# Vanishing Gradient Problem



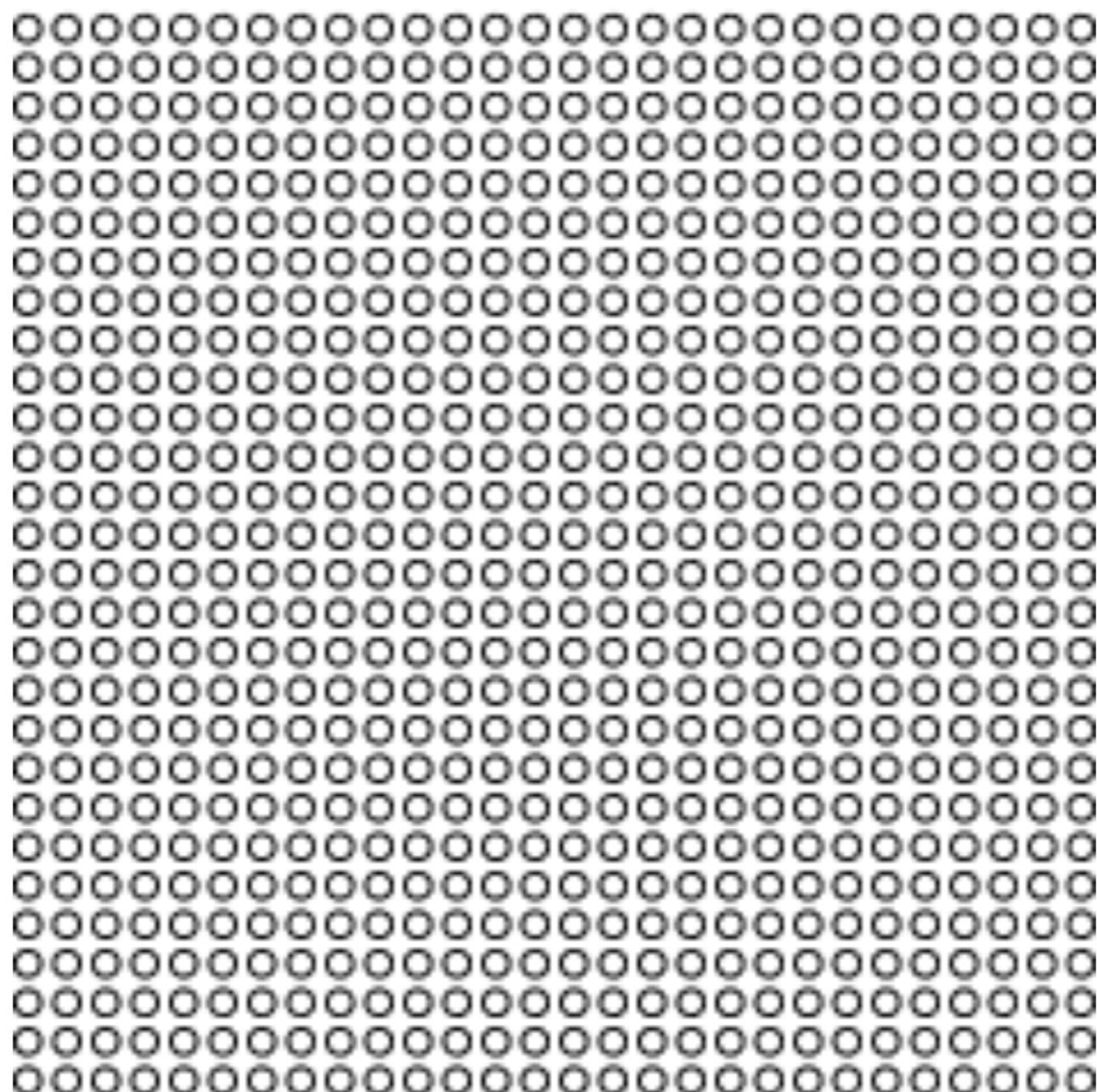
# Vanishing Gradient Problem



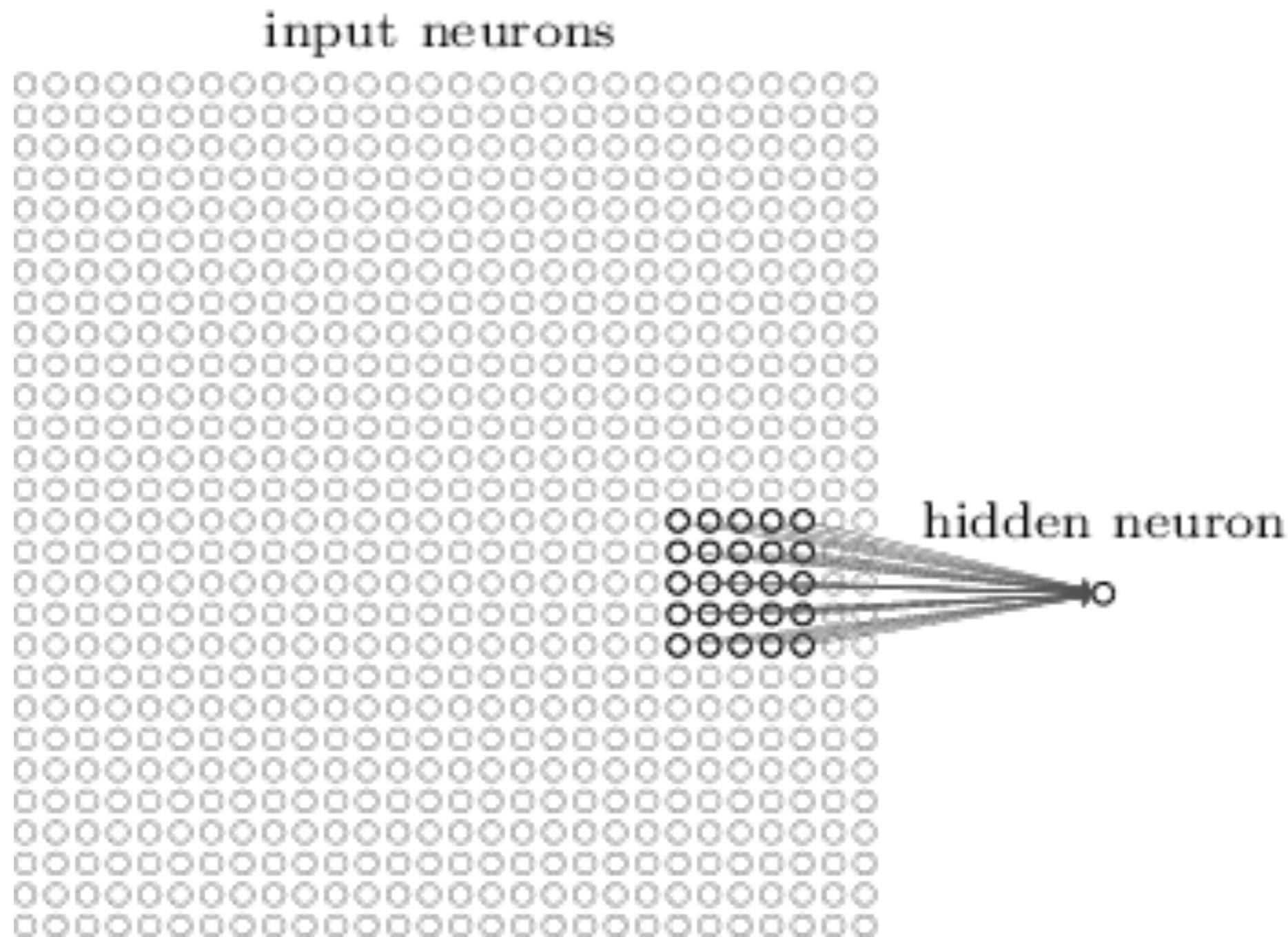
What now?

# Images as Input

input neurons



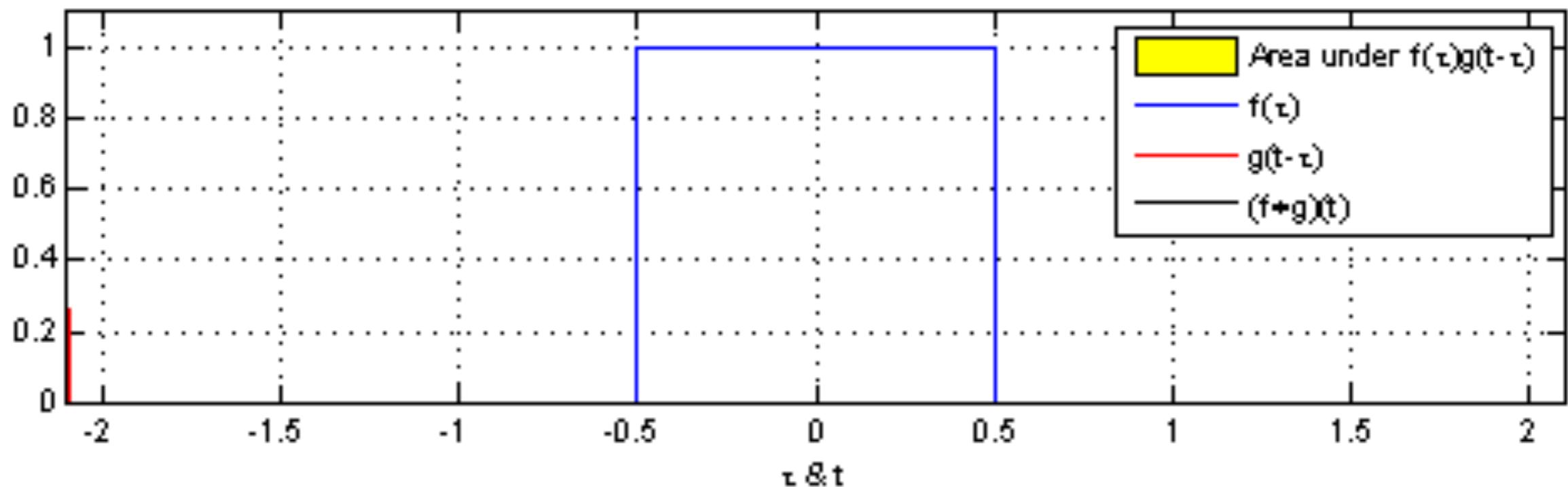
# Local receptive Fields



# Detour Convolution

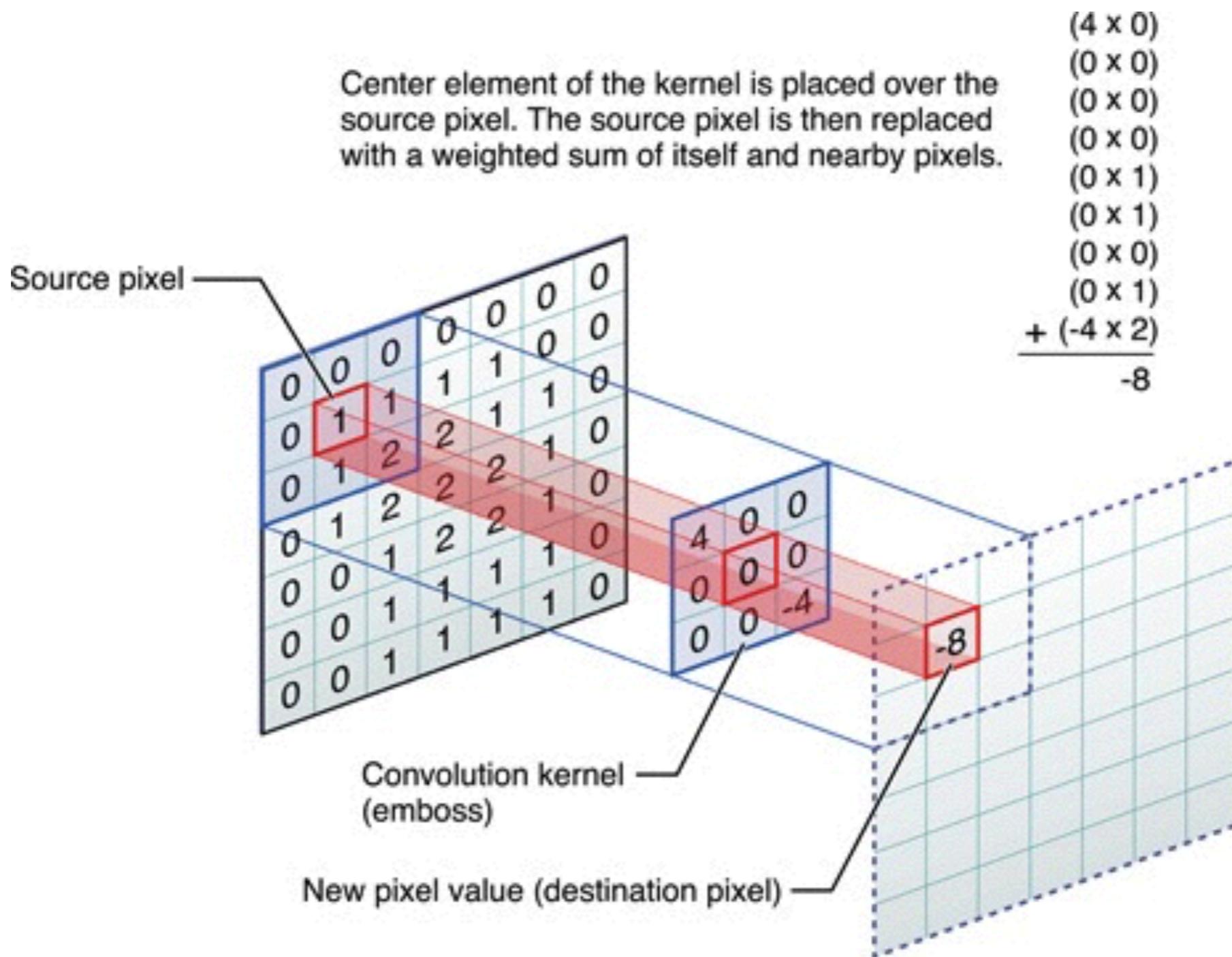
# Convolution

- is a mathematical operation of two functions
- and computes the integral
- of the point wise multiplication of the two functions.



# Convolution

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.



# Convolution

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

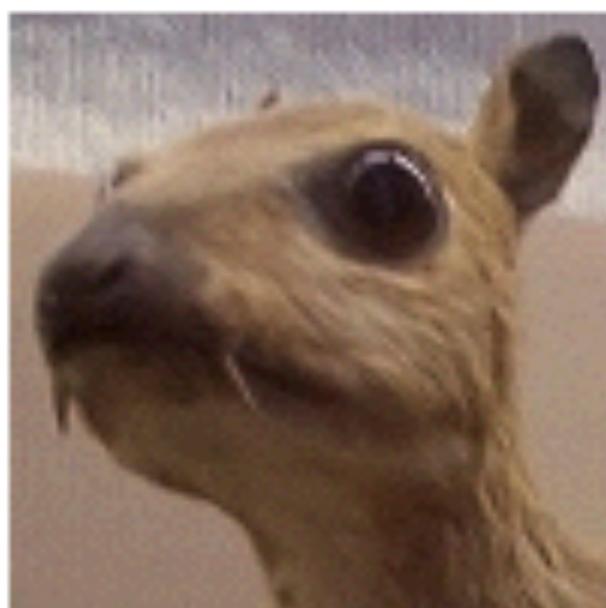
Image

4		

Convolved  
Feature

# Convolution

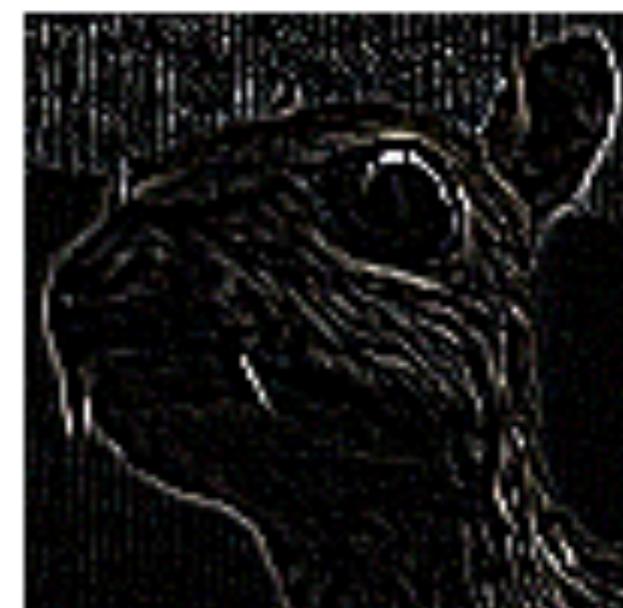
Input image



Convolution  
Kernel

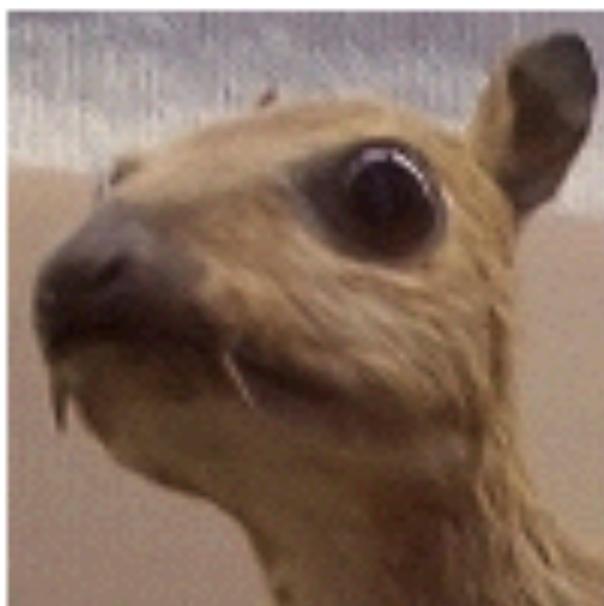
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Feature map



# Convolution

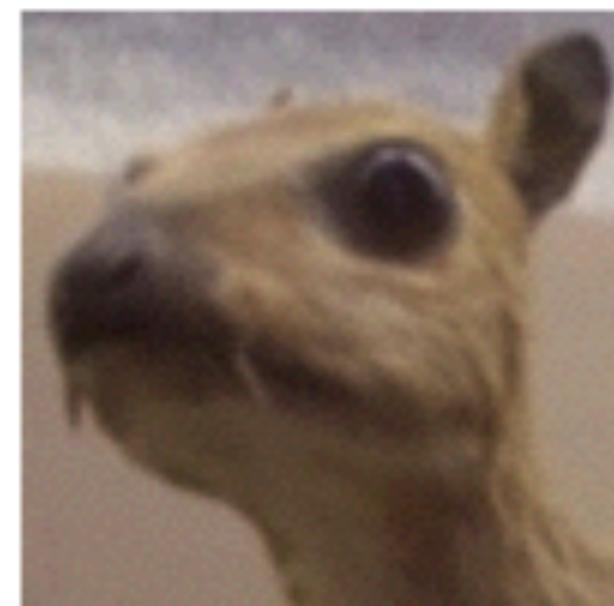
Input image



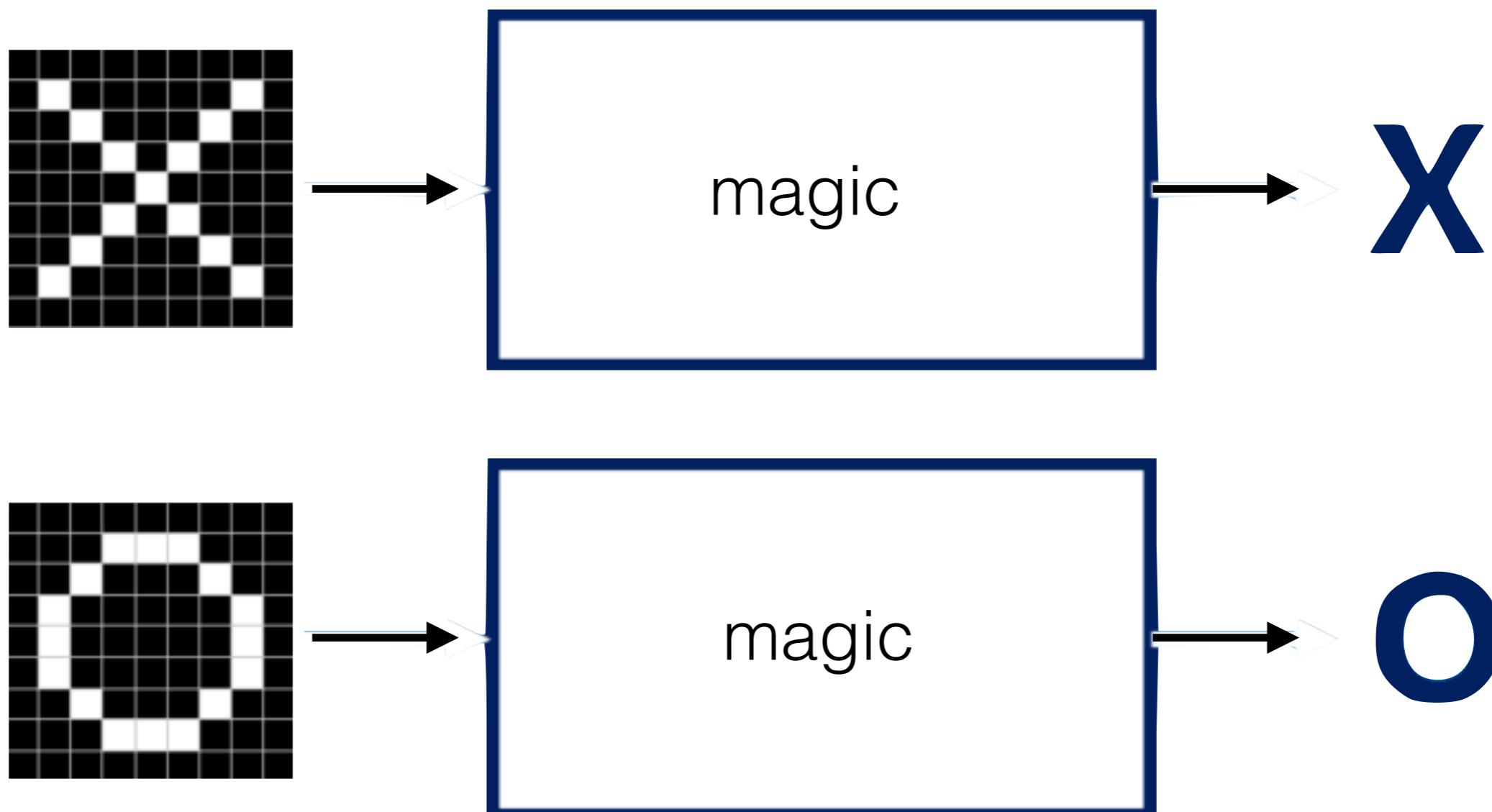
Kernel

$$\begin{pmatrix} \frac{1}{16} & \frac{1}{8} & \frac{1}{16} \\ \frac{1}{8} & \frac{1}{4} & \frac{1}{8} \\ \frac{1}{16} & \frac{1}{8} & \frac{1}{16} \end{pmatrix}$$

Feature map



# Object detection with convolutions



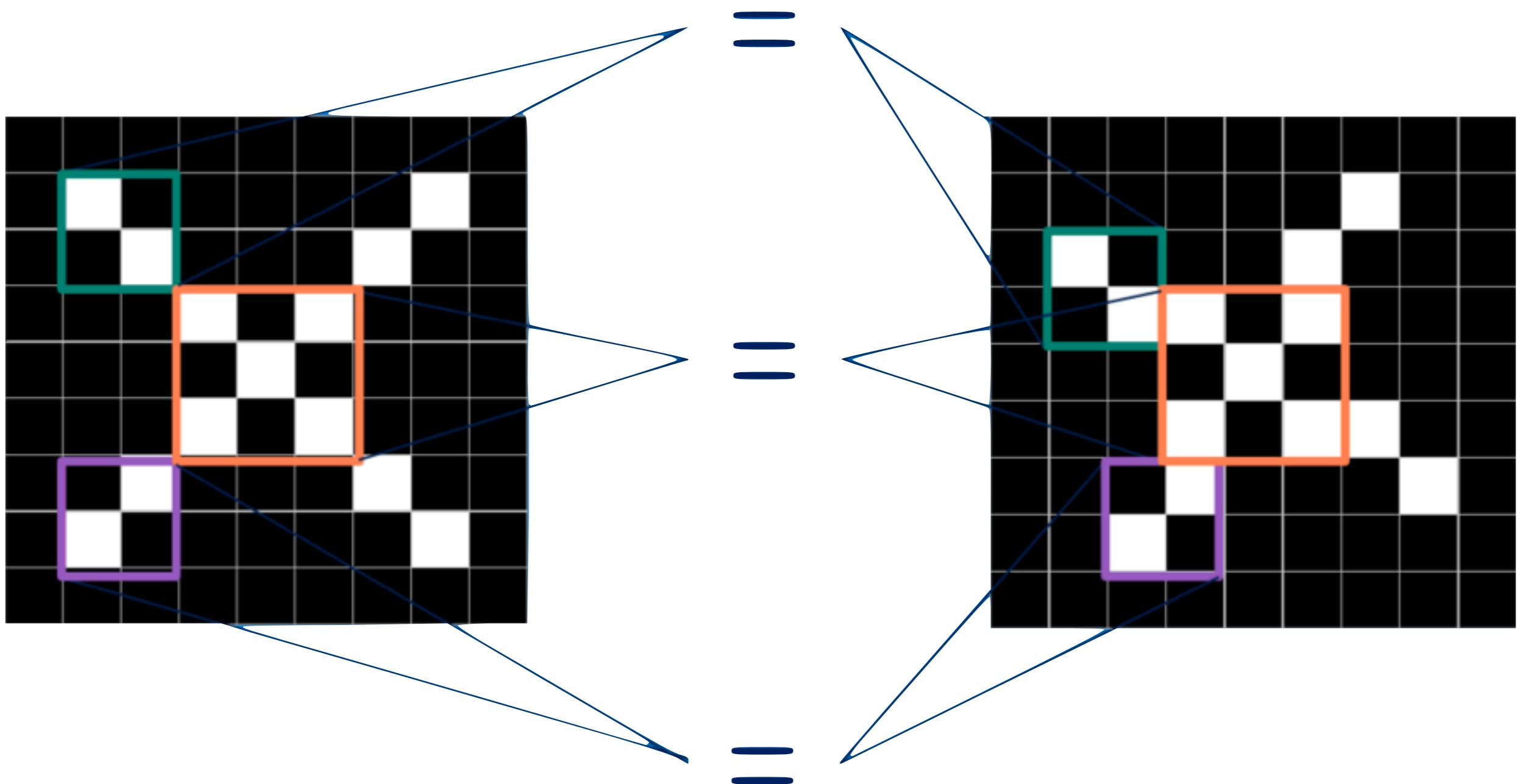
# Object detection with convolutions

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

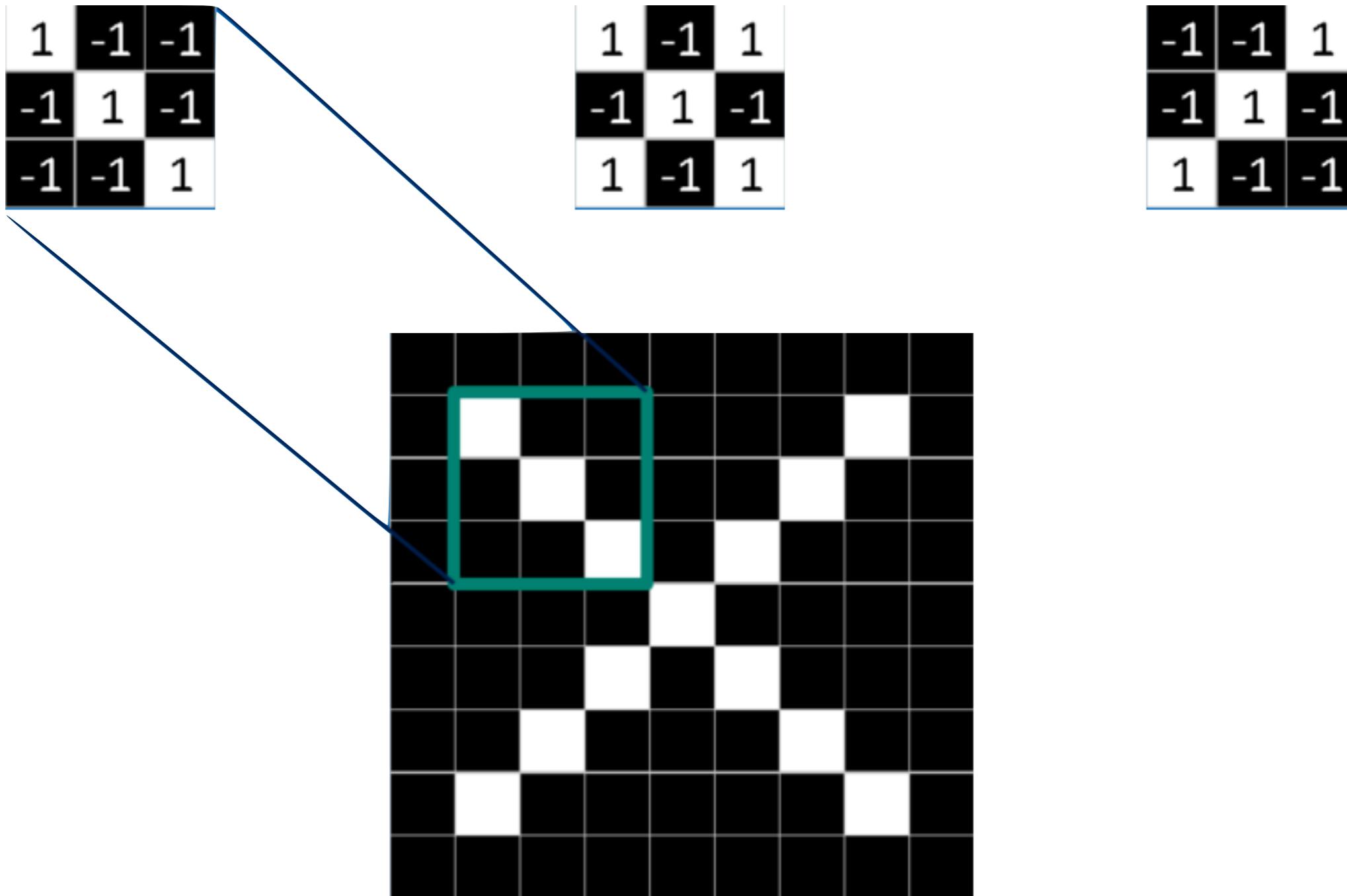


-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	-1	-1	-1	-1	-1	-1	1	-1
-1	1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	1	-1	1	-1	-1	-1
-1	-1	-1	1	1	-1	1	-1	-1
-1	-1	-1	-1	1	-1	1	-1	-1
-1	-1	-1	-1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	1	1	-1
-1	-1	-1	1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	-1	-1	-1

# Object detection with convolutions



# Object detection with convolutions



# Object detection with convolutions

1	-1	-1
-1	1	-1
-1	-1	1

filter

$$1 \times 1 = 1$$

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

1	1	1
1	1	1
1	1	1

result

# Object detection with convolutions

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	
-1	-1	1	-1	-1	-1	1	-1	-1	
-1	-1	-1	1	-1	1	-1	-1	-1	
-1	-1	-1	-1	1	-1	-1	-1	-1	
-1	-1	-1	-1	-1	1	-1	-1	-1	
-1	-1	-1	1	-1	1	-1	-1	-1	
-1	1	-1	-1	-1	-1	-1	1	-1	
-1	-1	-1	-1	-1	-1	-1	-1	-1	



1	-1	-1
-1	1	-1
-1	-1	1

=

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	
-1	-1	1	-1	-1	-1	1	-1	-1	
-1	-1	-1	1	-1	1	-1	-1	-1	
-1	-1	-1	-1	1	-1	-1	-1	-1	
-1	-1	-1	-1	-1	1	-1	-1	-1	
-1	-1	-1	1	-1	1	-1	-1	-1	
-1	1	-1	-1	-1	-1	-1	1	-1	
-1	-1	-1	-1	-1	-1	-1	-1	-1	



1	-1	1
-1	1	-1
1	-1	1

=

0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	
-1	-1	1	-1	-1	-1	1	-1	-1	
-1	-1	-1	1	-1	1	-1	-1	-1	
-1	-1	-1	-1	1	-1	-1	-1	-1	
-1	-1	-1	-1	-1	1	-1	-1	-1	
-1	-1	-1	1	-1	1	-1	-1	-1	
-1	1	-1	-1	-1	-1	-1	1	-1	
-1	-1	-1	-1	-1	-1	-1	-1	-1	

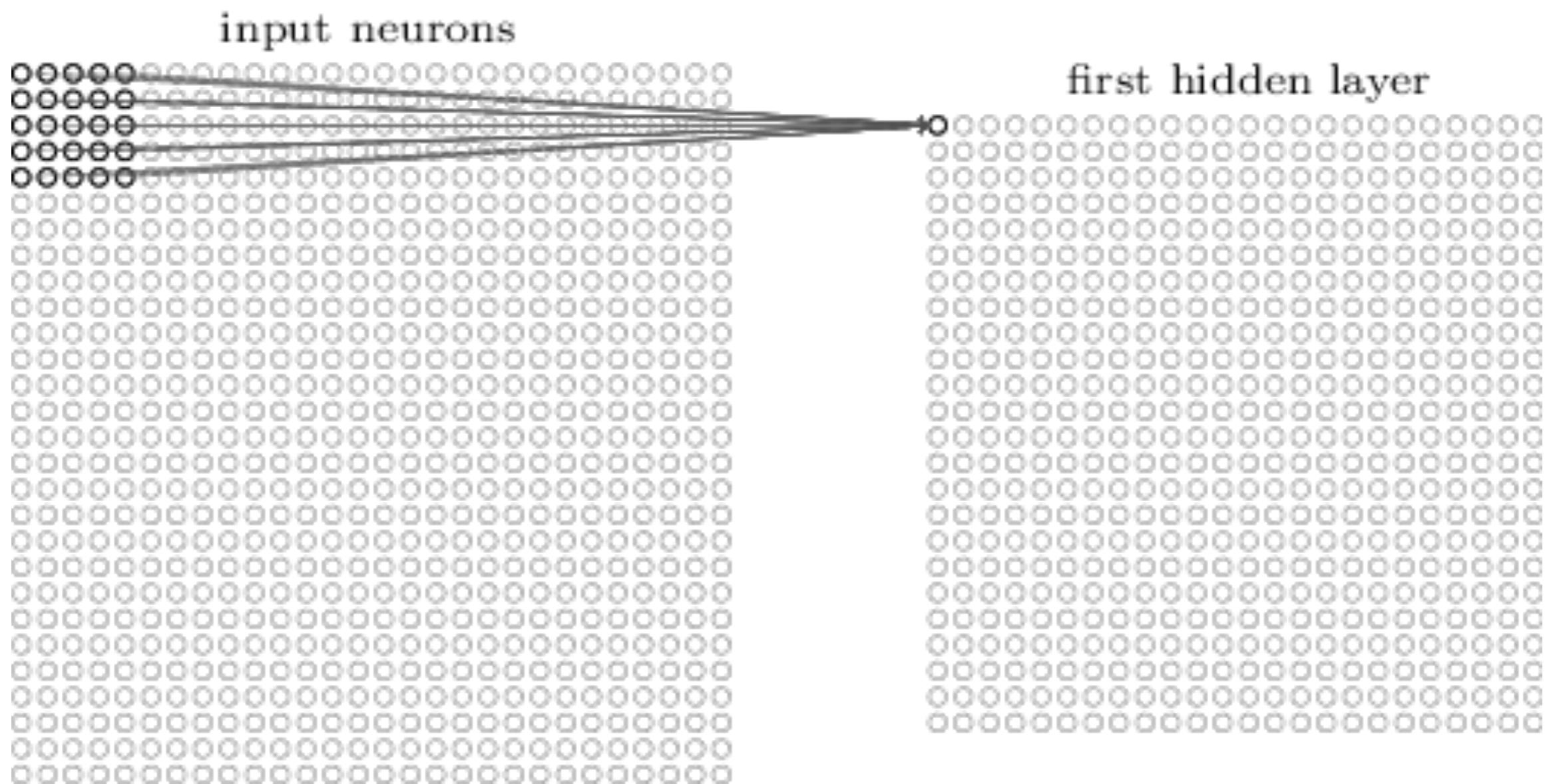


-1	-1	1
-1	1	-1
1	-1	-1

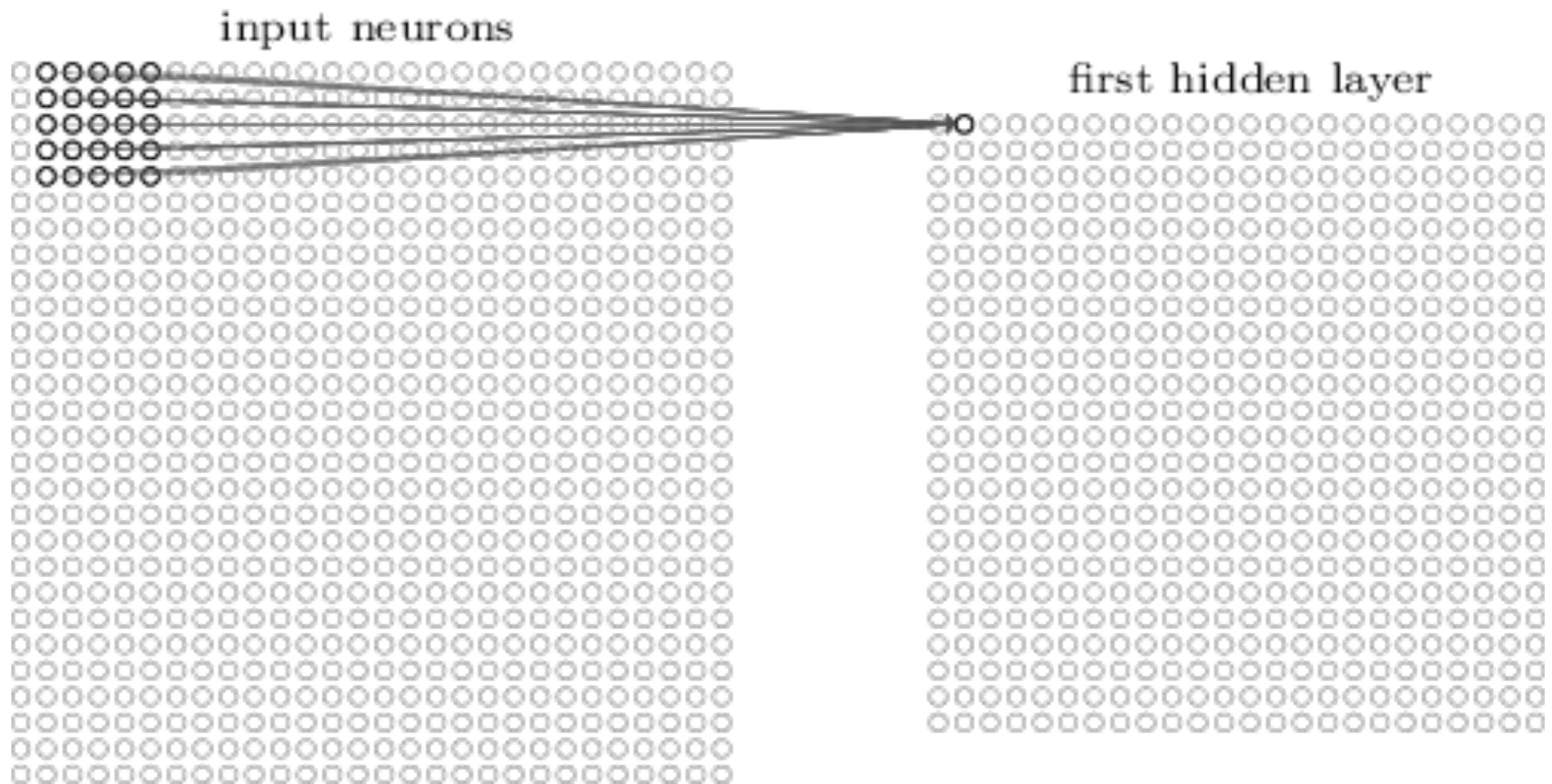
=

0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33

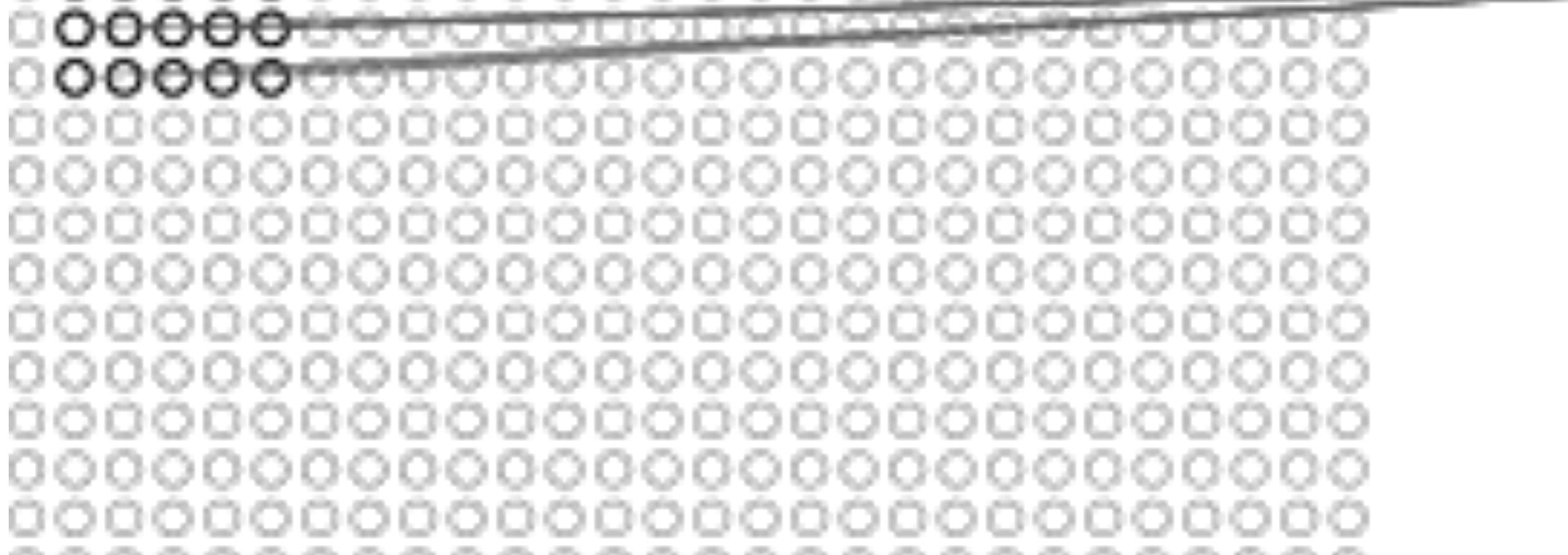
# Shared weights and biases



# Shared weights and biases



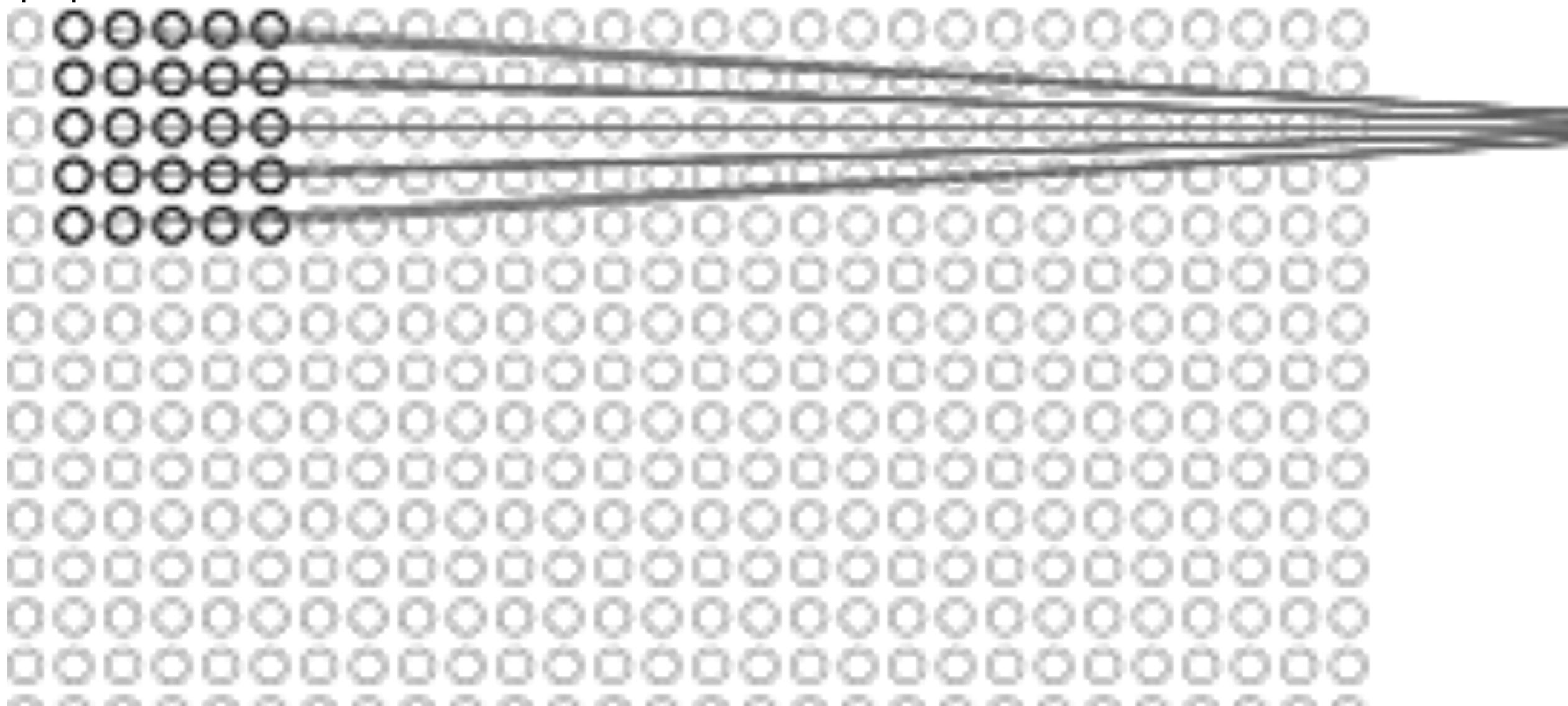
# Strides



# Strides

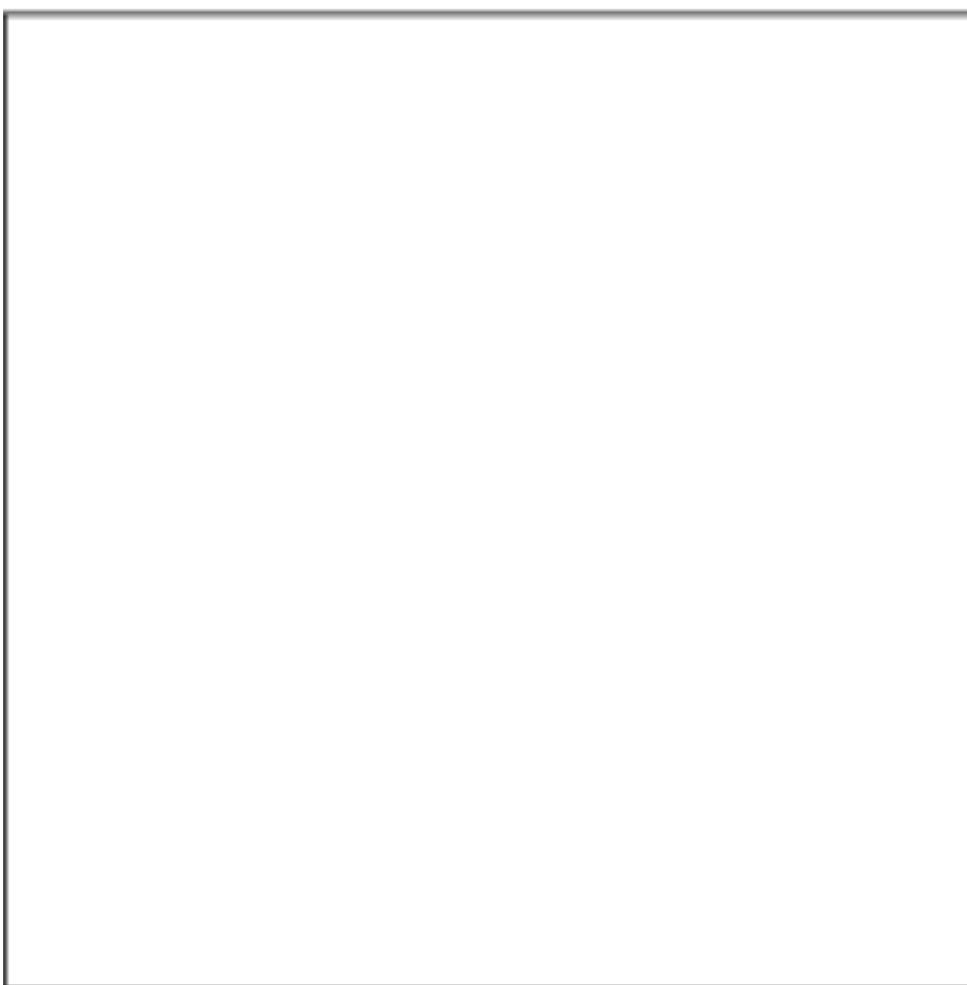
Stride of one

H

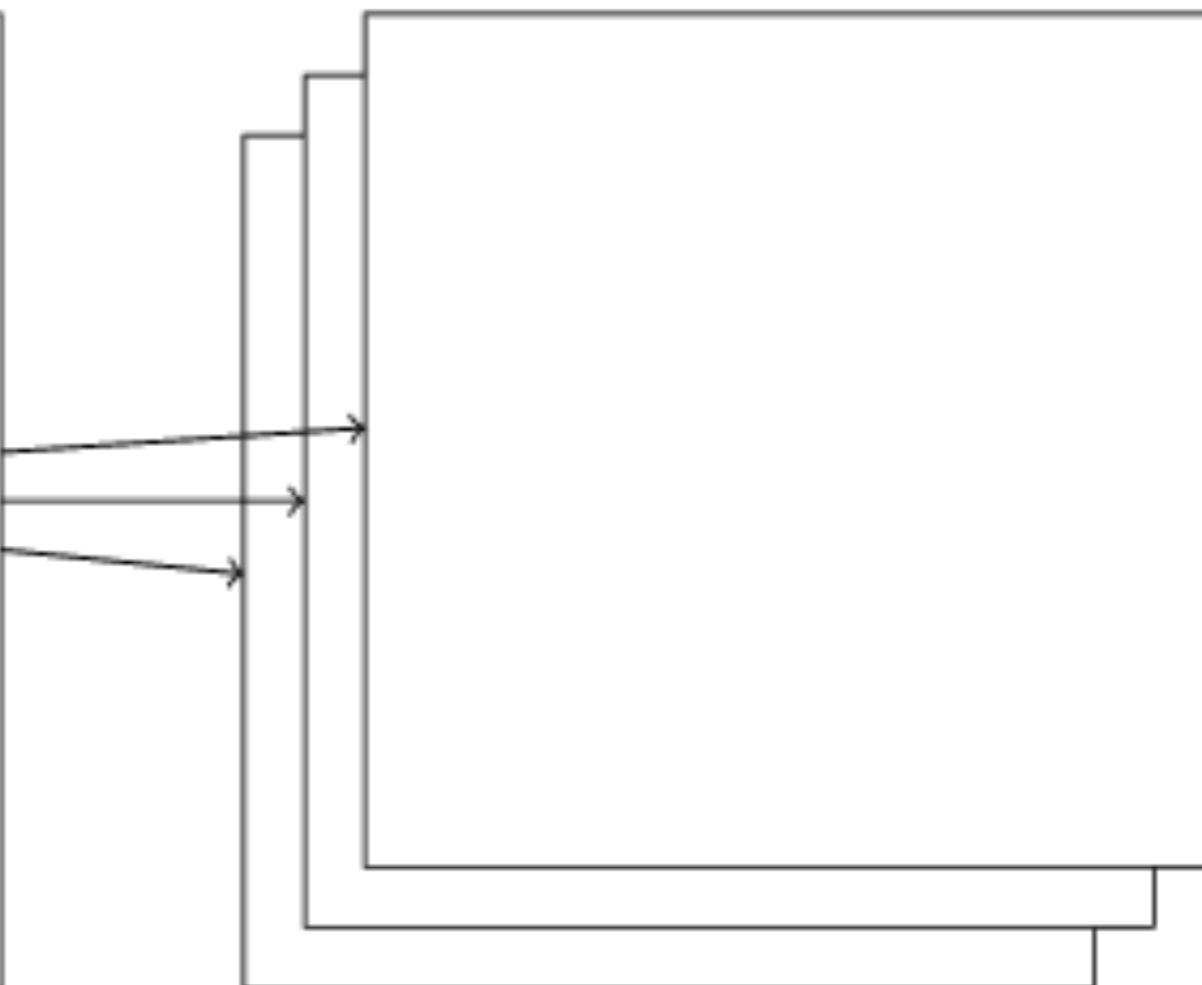


# Filter maps

$28 \times 28$  input neurons



first hidden layer:  $3 \times 24 \times 24$  neurons



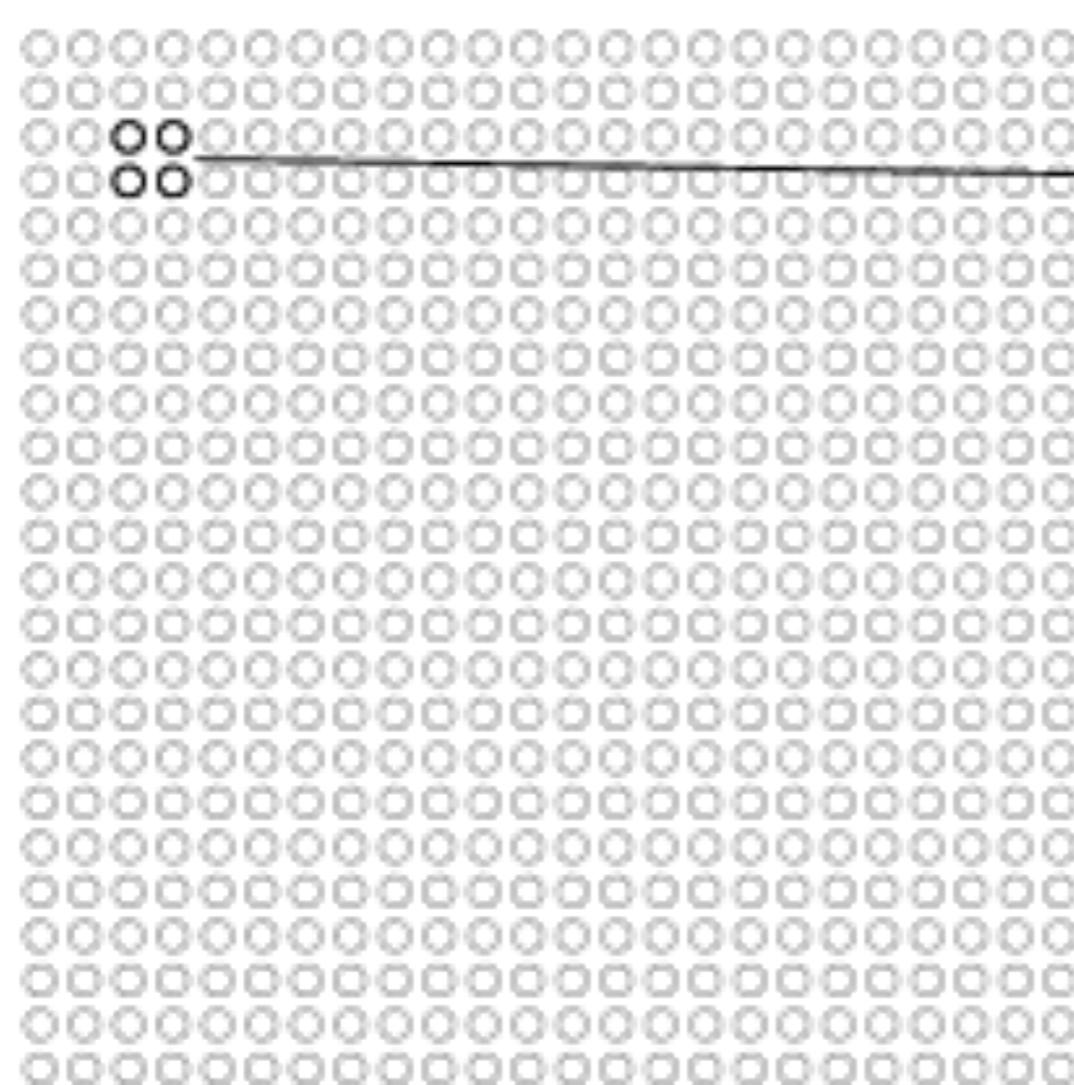
# Learned Filters



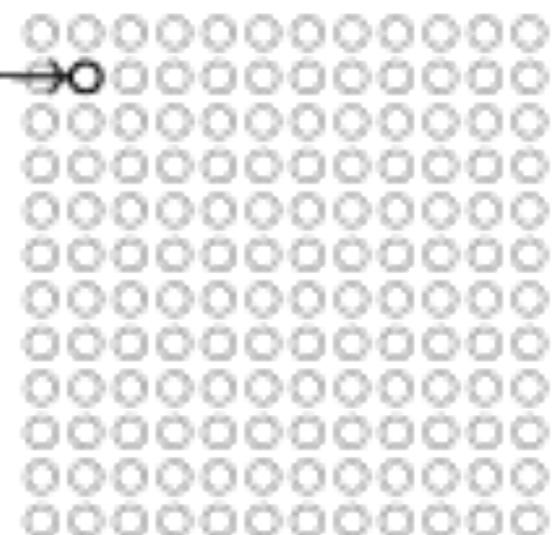
<http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks>

# Pooling

hidden neurons (output from feature map)

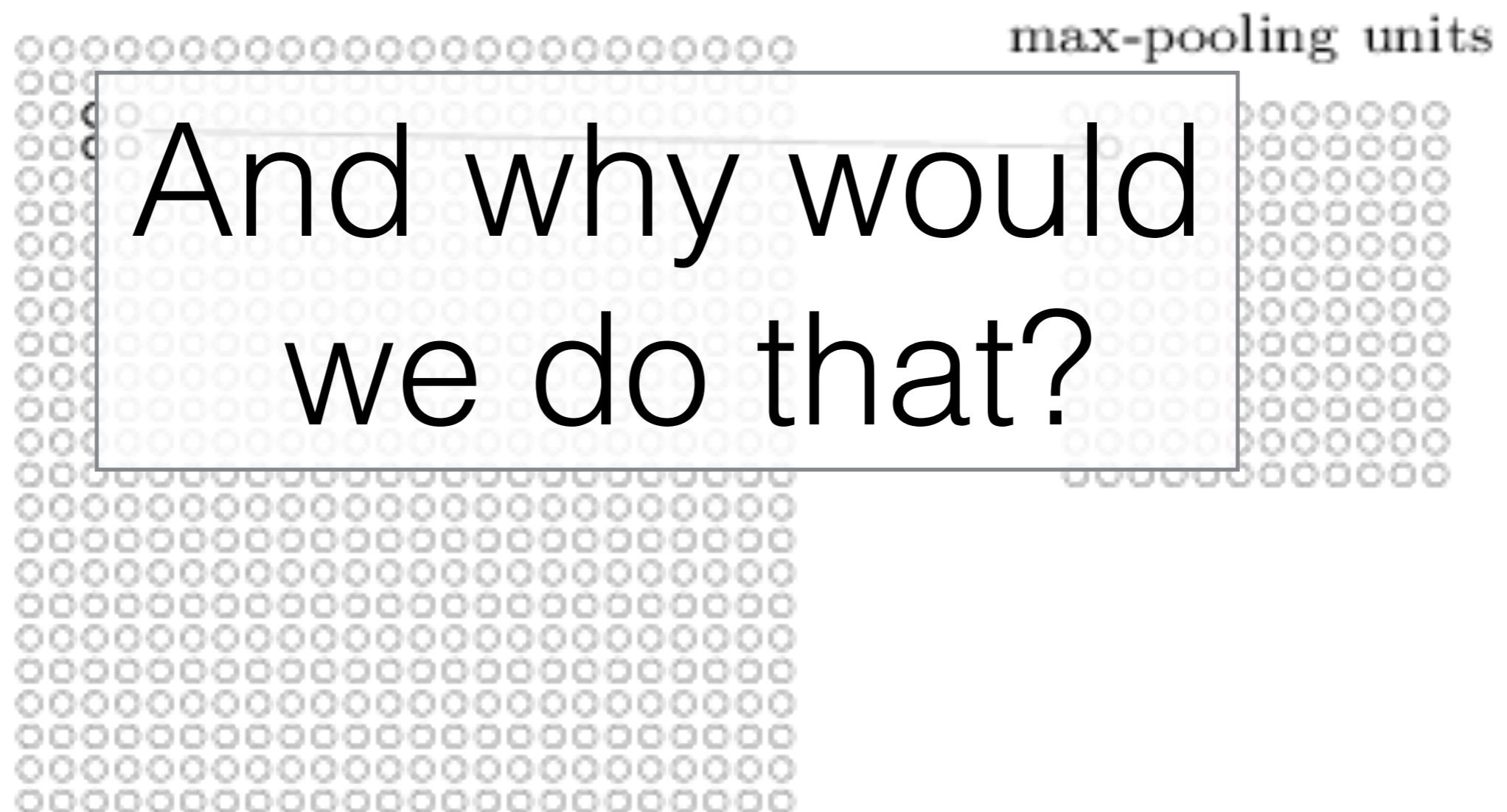


max-pooling units



# Pooling

hidden neurons (output from feature map)



# Pooling

hidden neurons (output from feature map)

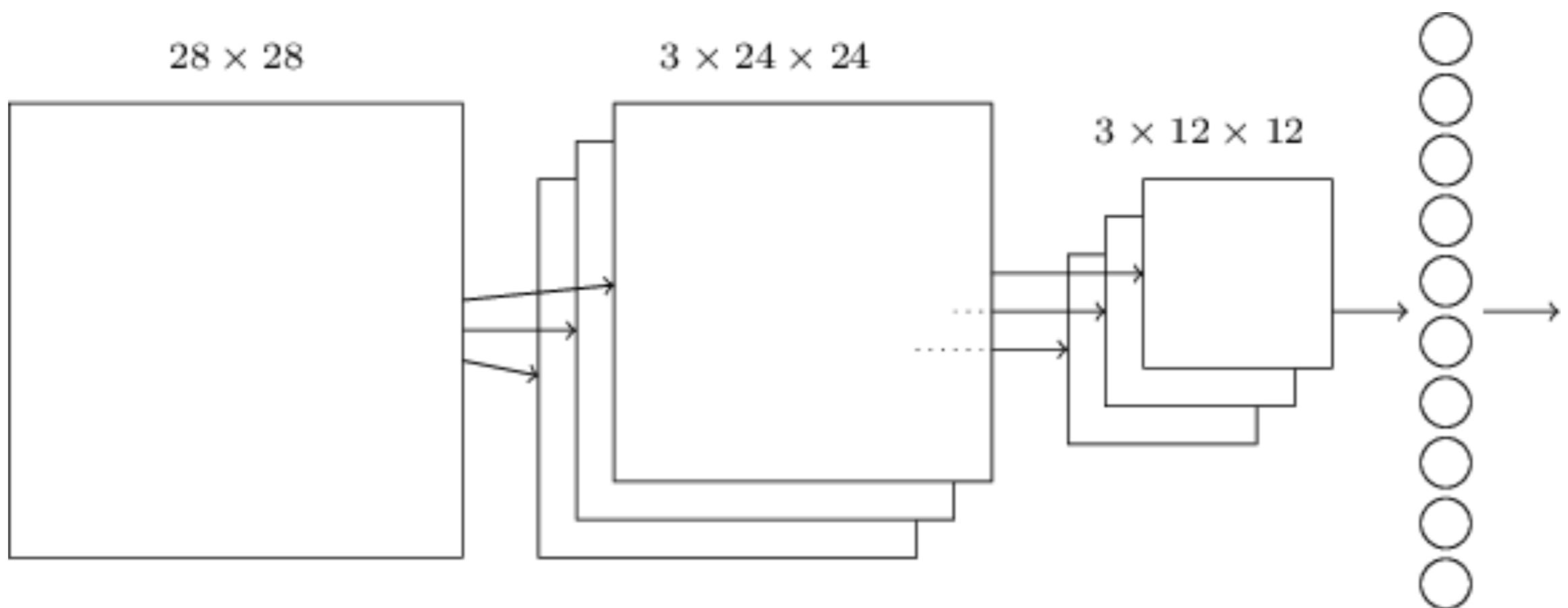
oooooooooooooooooooo

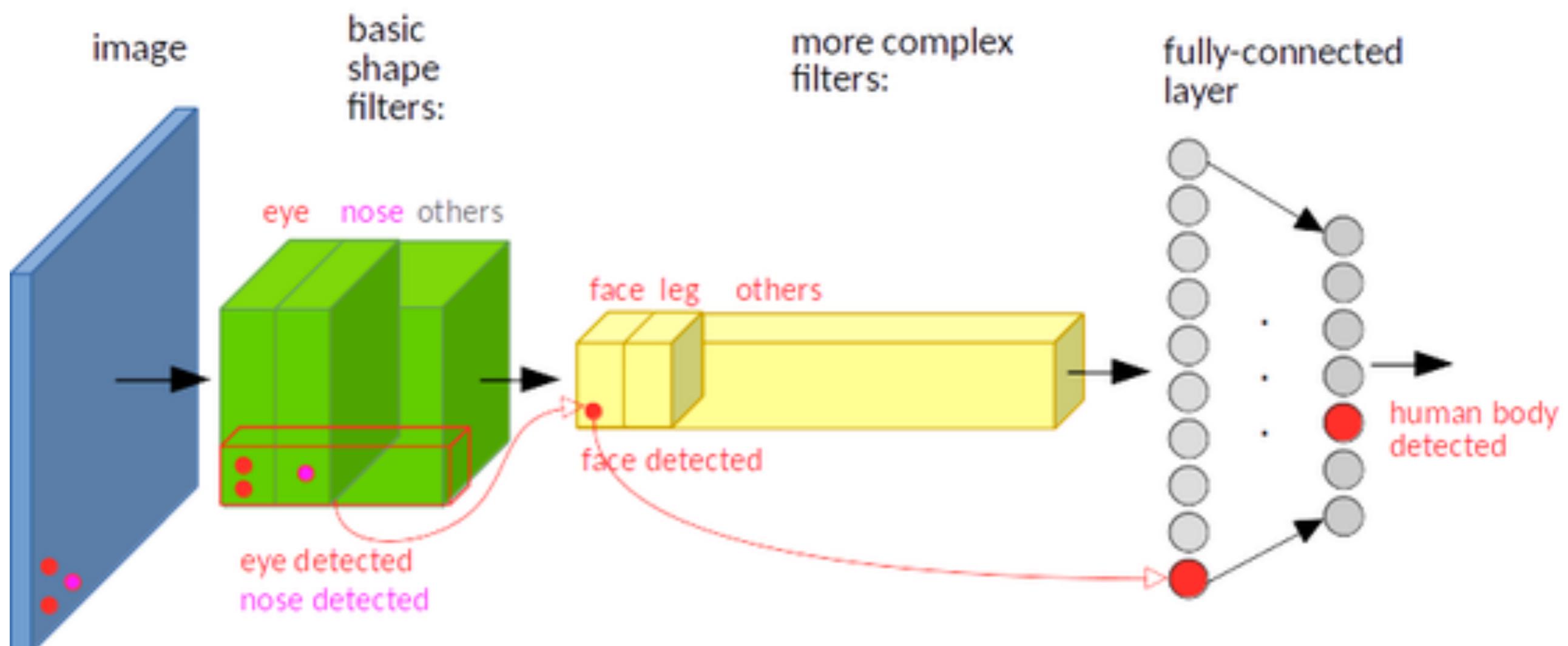
max-pooling units

And why would  
we do that?

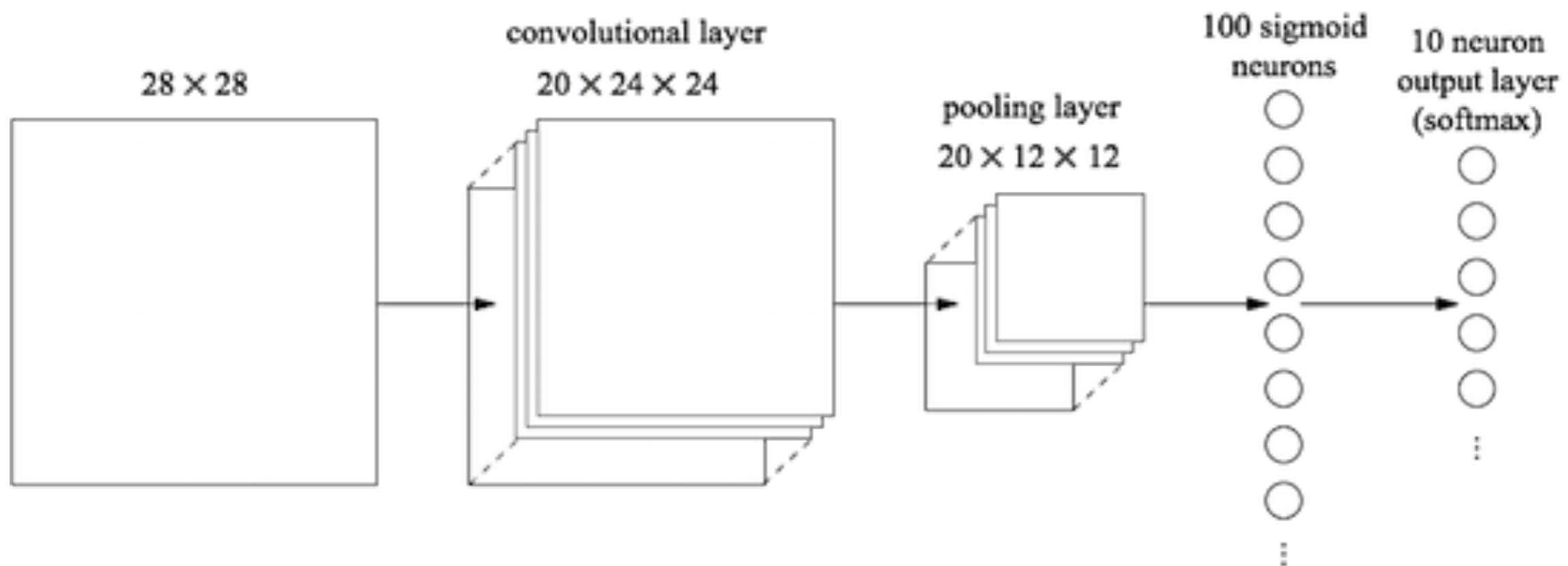
=> local invariance &  
reduced parameters

# Convolutional Neural Networks

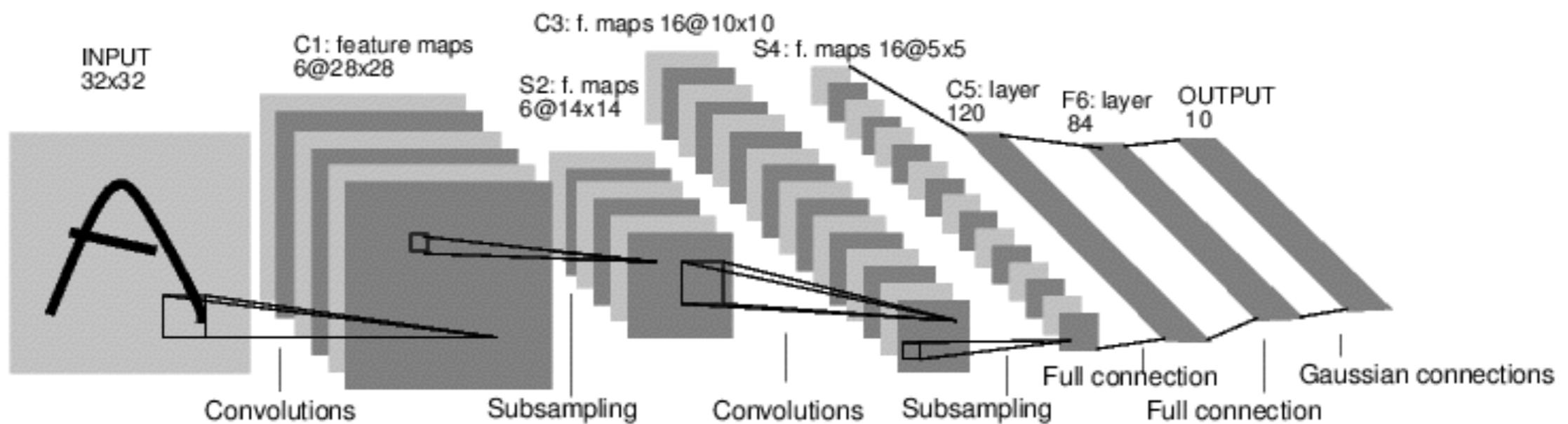




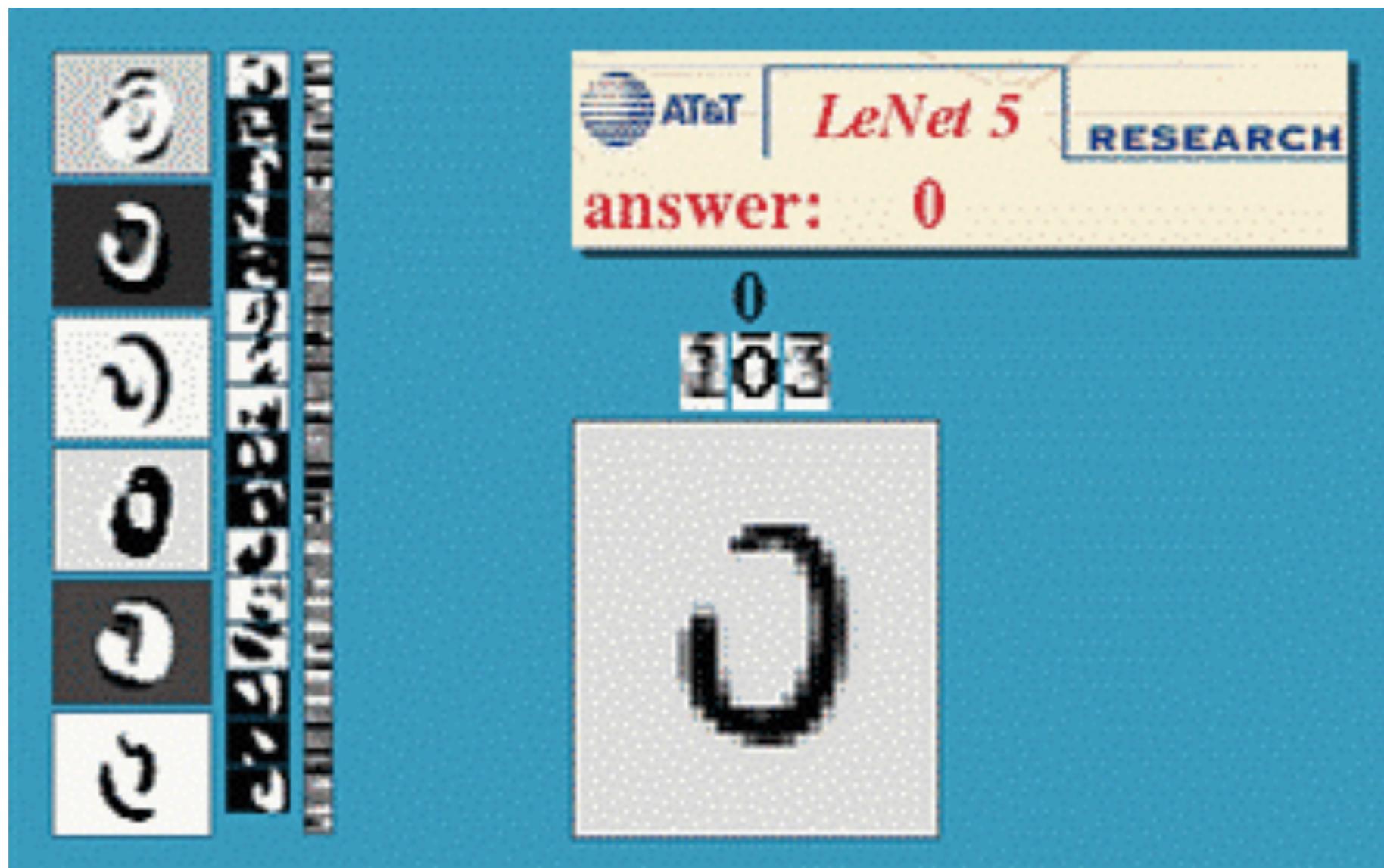
# Convolutional Neural Networks



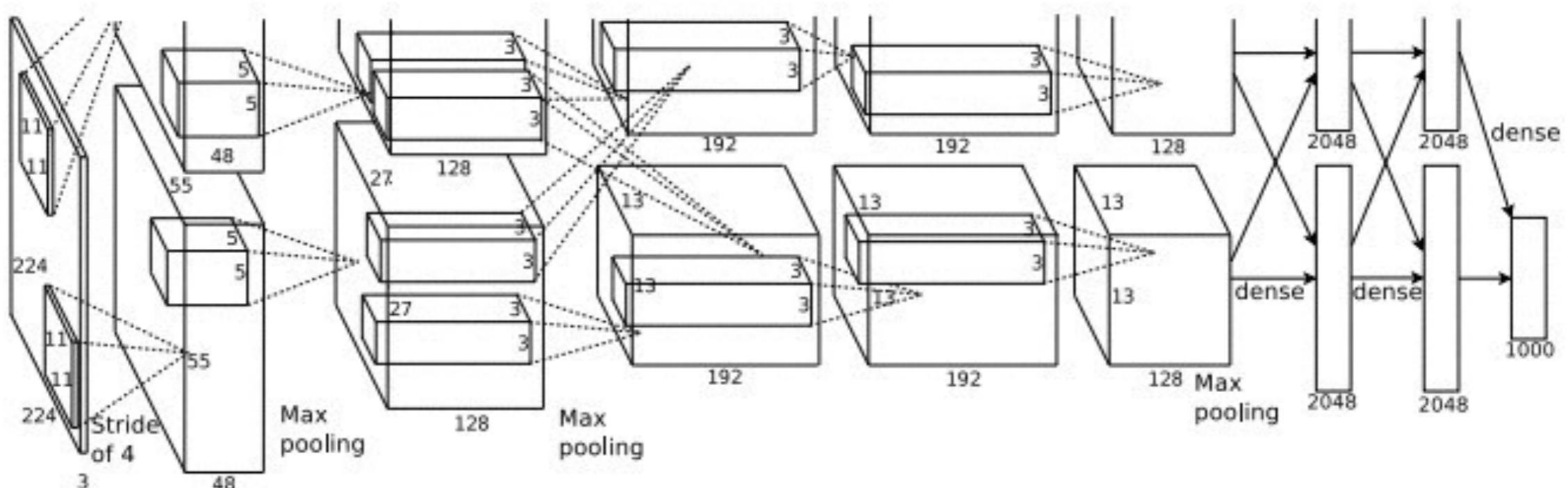
# LeNet5 1998



# LeNet5 1998

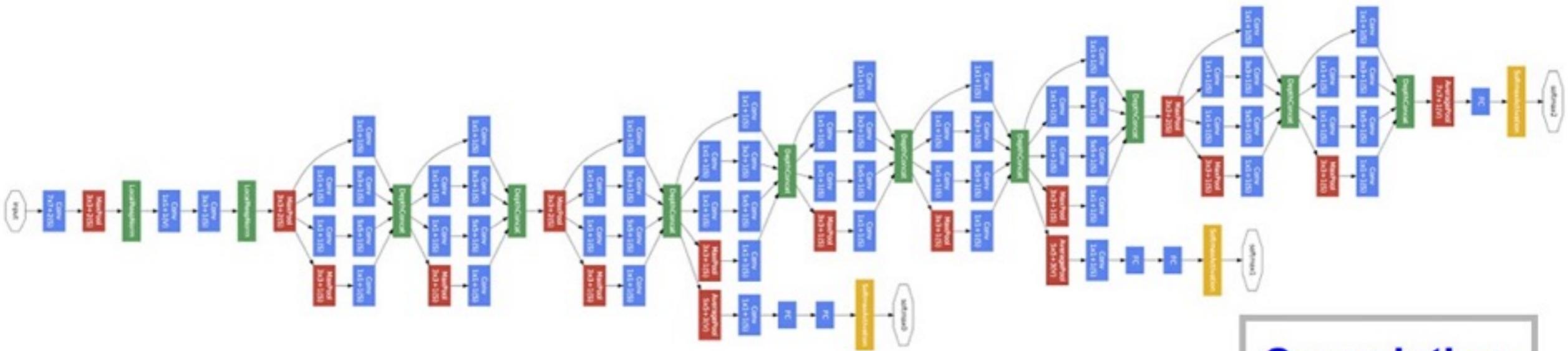


# Alexnet 2012



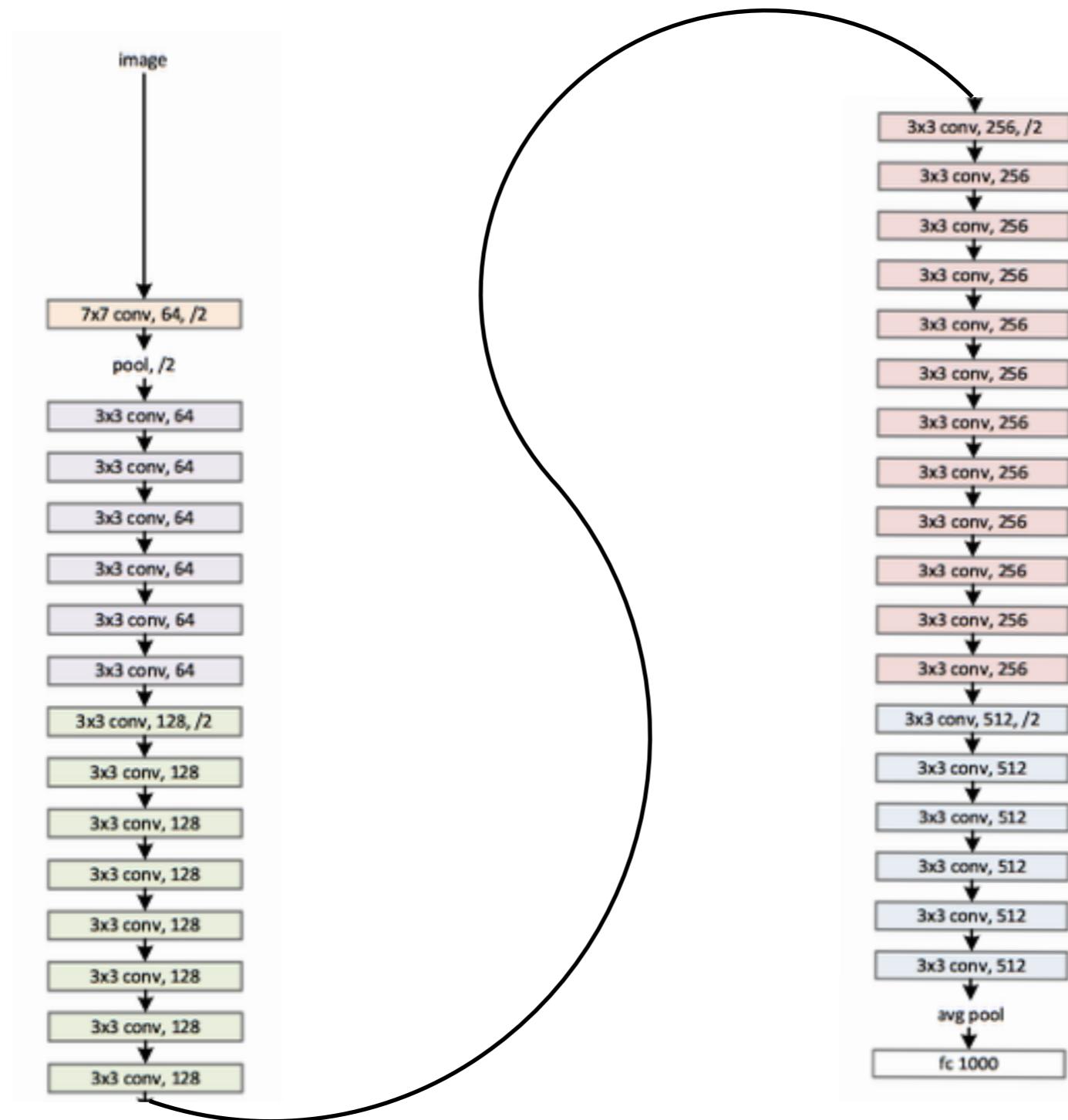
<http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks>

# Google LeNet 2014

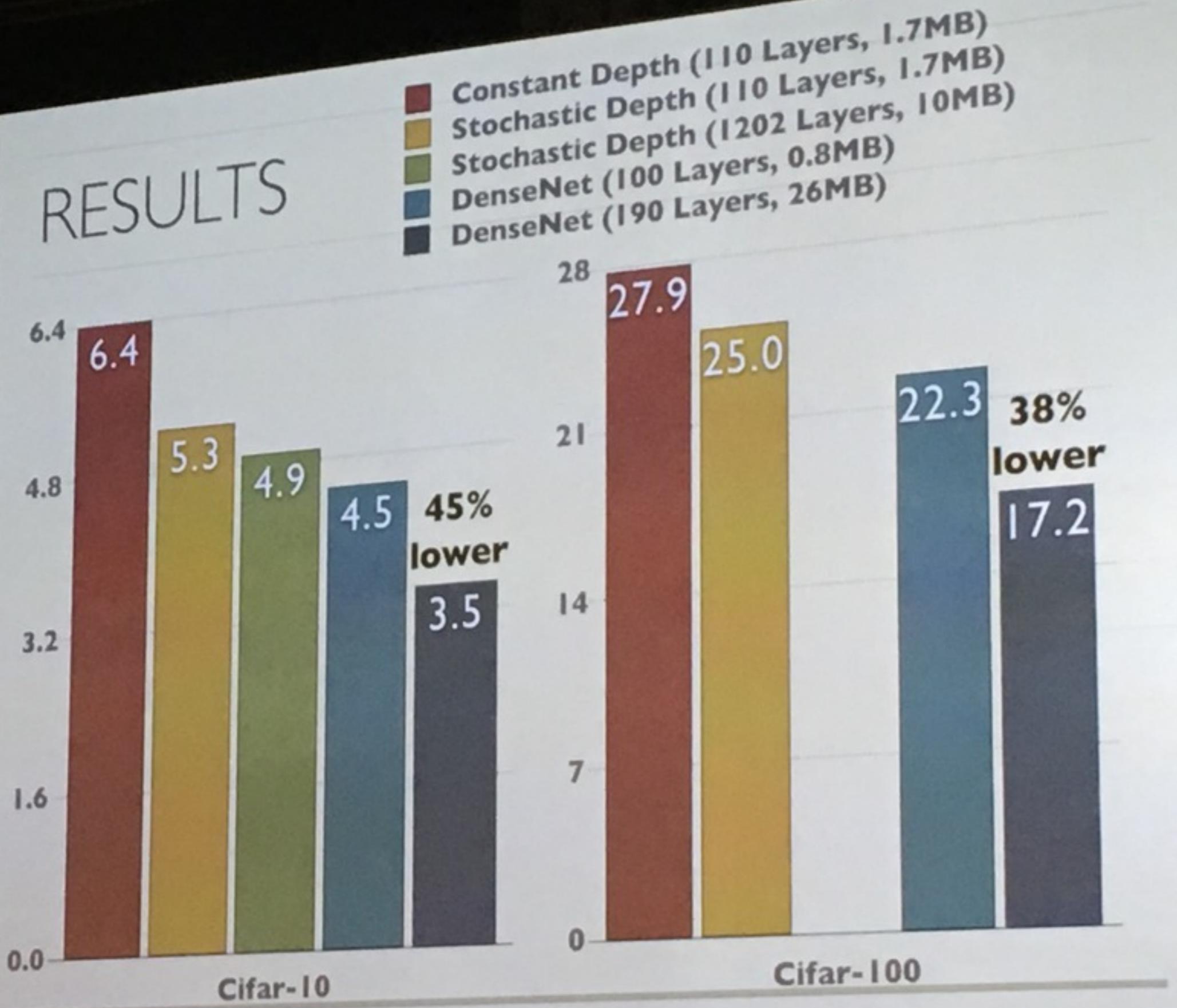


**Convolution**  
**Pooling**  
**Softmax**  
**Other**

# Microsofts Deep Residual Networks 2015



# RESULTS



# RESULTS



# Frameworks

# `Hello World` in CUDA

```
#include <stdio.h>
const int N = 16;
const int blocksize = 16;
__global__
void hello(char *a, int *b)
{
    a[threadIdx.x] += b[threadIdx.x];
}
int main()
{
    char a[N] = "Hello \0\0\0\0\0\0";
    int b[N] = {15, 10, 6, 0, -11, 1,
0, 0, 0, 0, 0, 0, 0, 0};
    char *ad;
    int *bd;
    const int csize = N*sizeof(char);
    const int isize = N*sizeof(int);
    printf("%s", a);
    cudaMalloc( (void**)&ad, csize );
    cudaMalloc( (void**)&bd, isize );
    cudaMemcpy( ad, a, csize,
cudaMemcpyHostToDevice );
    cudaMemcpy( bd, b, isize,
cudaMemcpyHostToDevice );

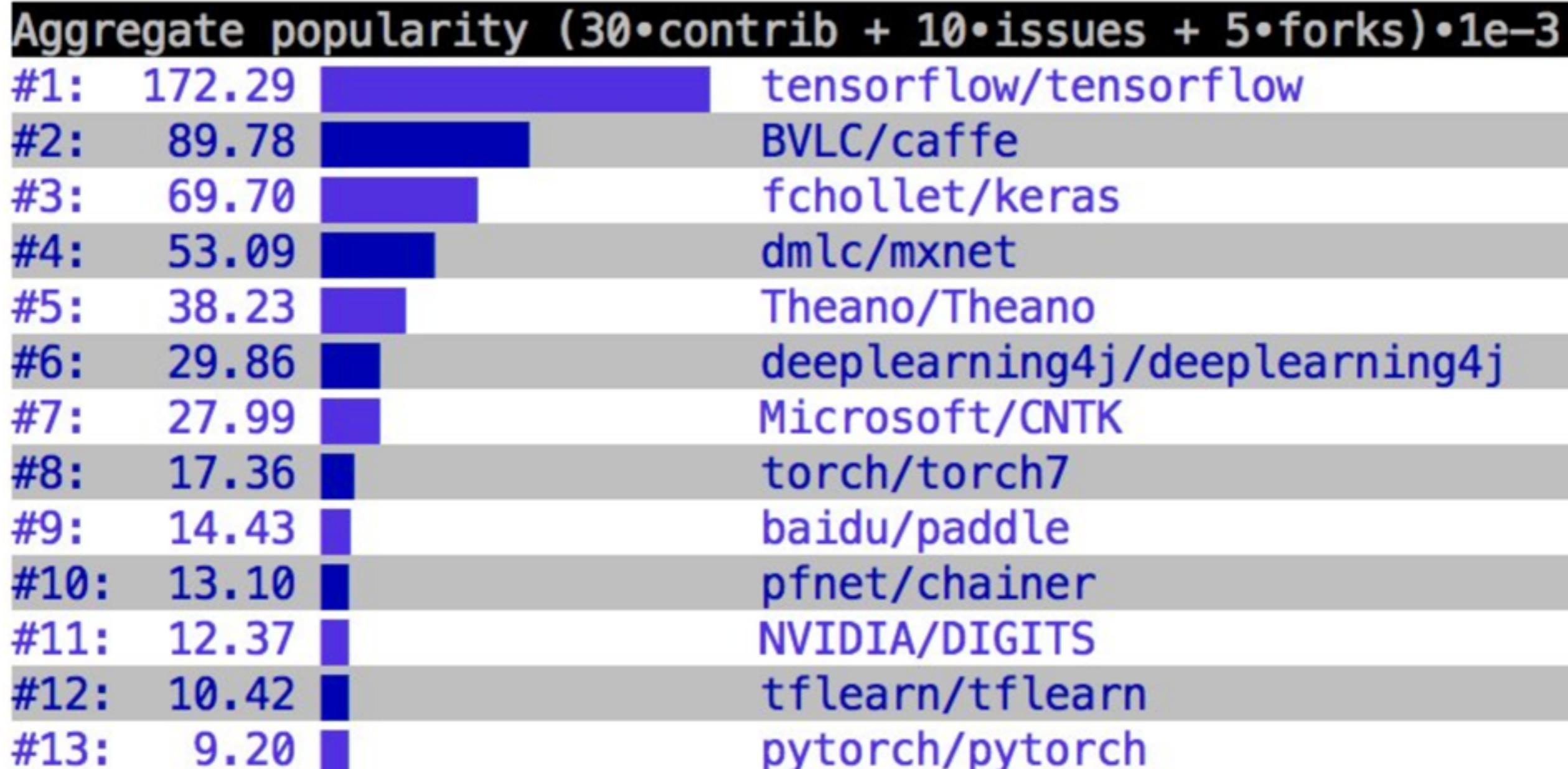
    dim3 dimBlock( blocksize, 1 );
    dim3 dimGrid( 1, 1 );
    hello<<<dimGrid,
dimBlock>>>(ad, bd);
    cudaMemcpy( a, ad, csize,
cudaMemcpyDeviceToHost );
    cudaFree( ad );
    cudaFree( bd );

    printf("%s\n", a);
    return EXIT_SUCCESS;
}
```

# How to decide what Deep Learning library to use?

- Features
  - `speciality` layers
  - most recent loss and activation functions
  - training ‘hacks’ (dropout)
  - helper functions
- Abstraction
- Speed
  - CPU/GPU speed
  - multi GPU support
- Ease of customization
- Base language
  - Python, C++, Java, Rust(!)...

# Deep learning libraries: Accumulated GitHub metrics

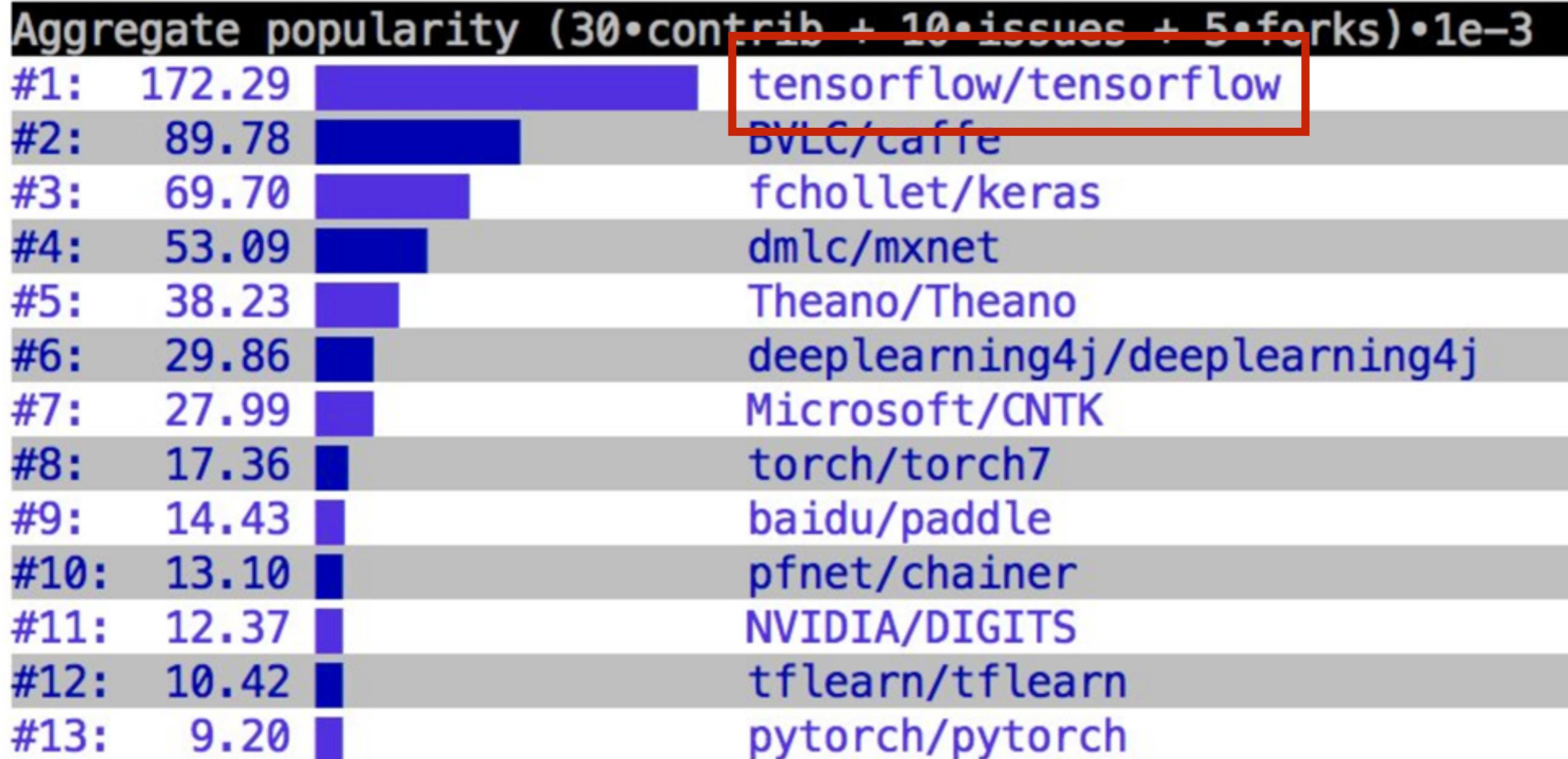


François Chollet @fchollet · Feb 11

Time for an update: what does the deep learning library landscape look like, seen from GitHub?

[pic.twitter.com/QDZyvLrYBd](https://pic.twitter.com/QDZyvLrYBd)

# Deep learning libraries: Accumulated GitHub metrics



François Chollet @fchollet · Feb 11

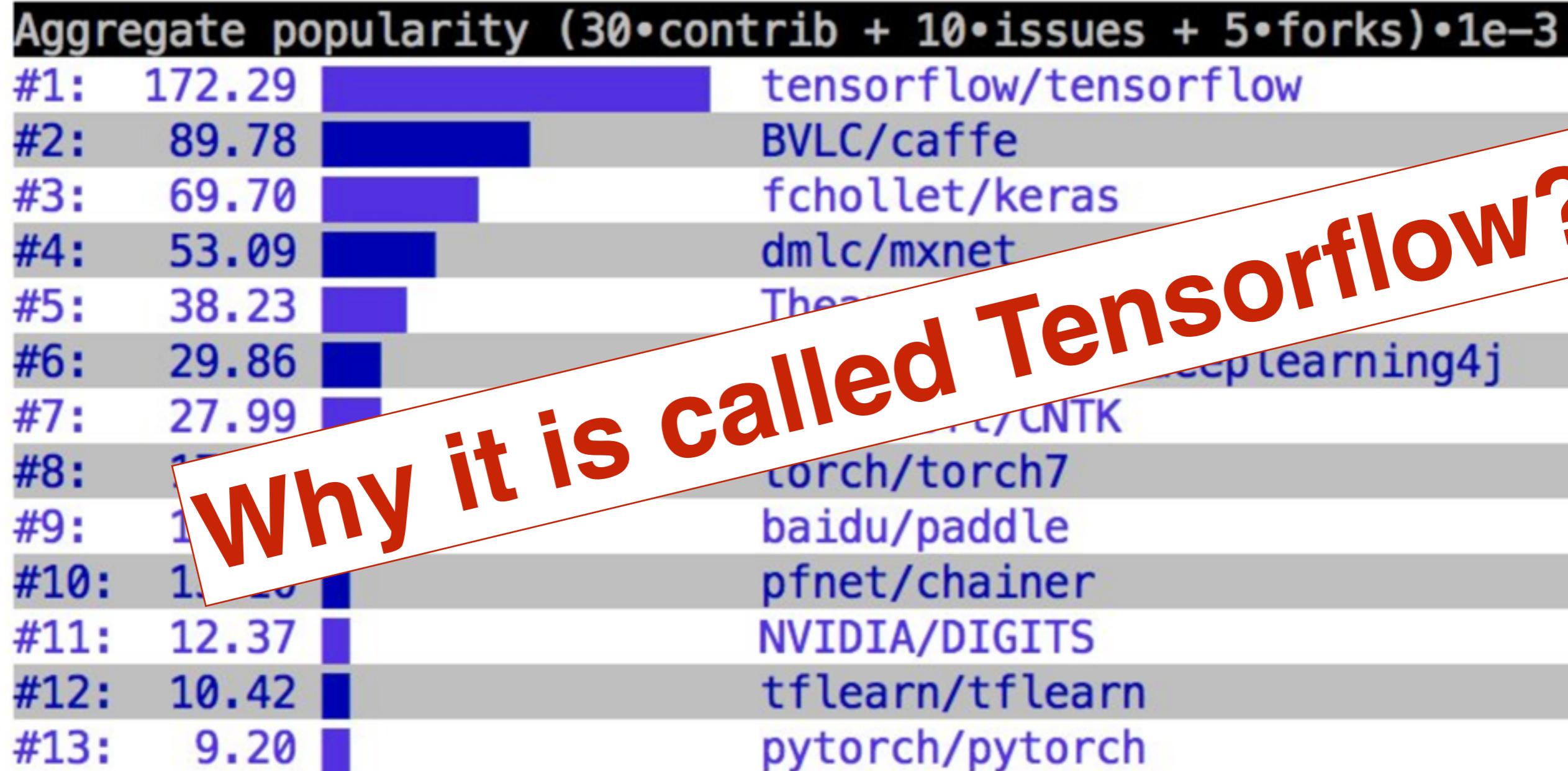
Time for an update: what does the deep learning library landscape look like, seen from GitHub?

[pic.twitter.com/QDZyvLrYBd](https://pic.twitter.com/QDZyvLrYBd)

# Tensor flow

- Features
  - `speciality` layers
  - most recent loss and activation functions 
  - training ‘hacks’ (dropout) 
  - helper functions 
- Abstraction
- Speed
  - CPU/GPU speed
  - multi GPU support 
- Ease of customization 
- Base language
  - Python , C++, Java, Rust(!)...

# Deep learning libraries: Accumulated GitHub metrics



François Chollet @fchollet · Feb 11

Time for an update: what does the deep learning library landscape look like, seen from GitHub?

[pic.twitter.com/QDZyvLrYBd](https://pic.twitter.com/QDZyvLrYBd)

# Init

```
import tensorflow as tf  
sess = tf.InteractiveSession()
```

# Setup variables

```
x = tf.placeholder(tf.float32,  
                  shape=[None, 784])  
y_ = tf.placeholder(tf.float32,  
                  shape=[None, 10])  
W = tf.Variable(tf.zeros([784, 10]))  
b = tf.Variable(tf.zeros([10]))  
  
sess.run(tf.global_variables_initializer())
```

# Define the network

```
y = tf.matmul(x, W) + b
softmax = tf.nn.\
    softmax_cross_entropy_with_logits(y_, y)
cross_entropy = tf.reduce_mean(softmax)
GD = tf.train.GradientDescentOptimizer
train_step = GD().minimize(cross_entropy)
```

# Train

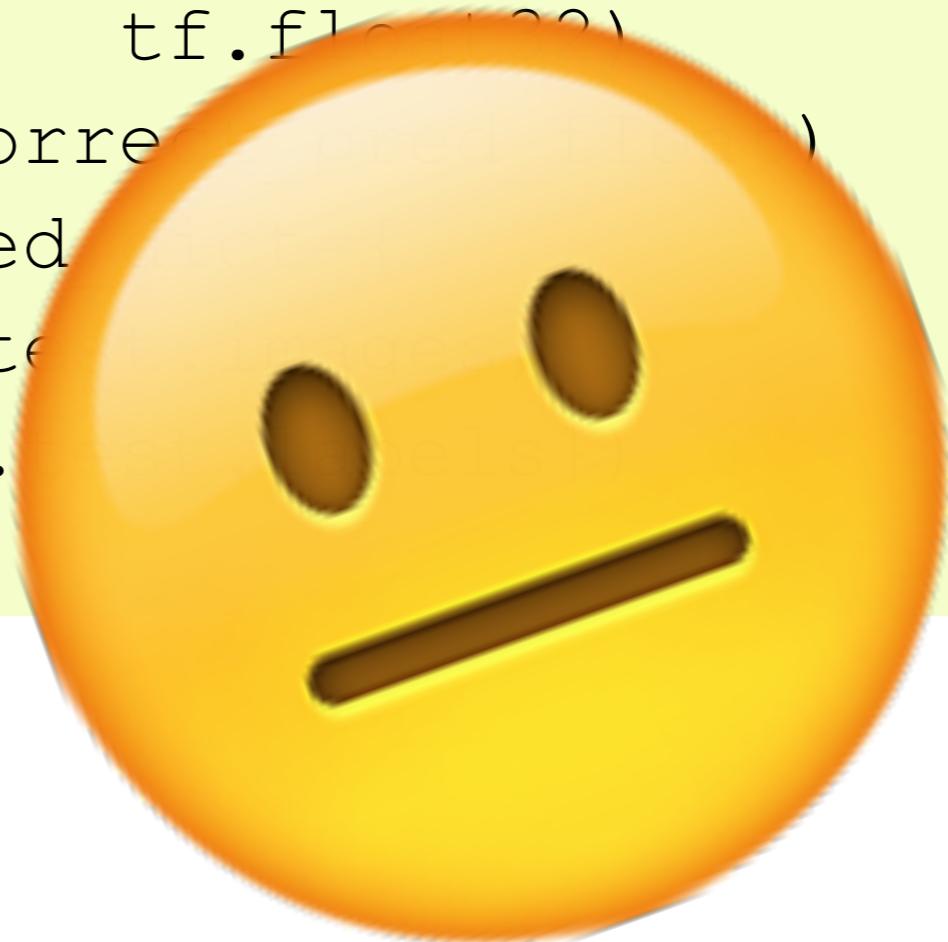
```
GD = tf.train.GradientDescentOptimizer  
train_step = GD().minimize(cross_entropy)  
for _ in range(1000):  
    batch = mnist.train.next_batch(100)  
    train_step.run(feed_dict={x: batch[0],  
                            y_: batch[1]})
```

# Evaluate

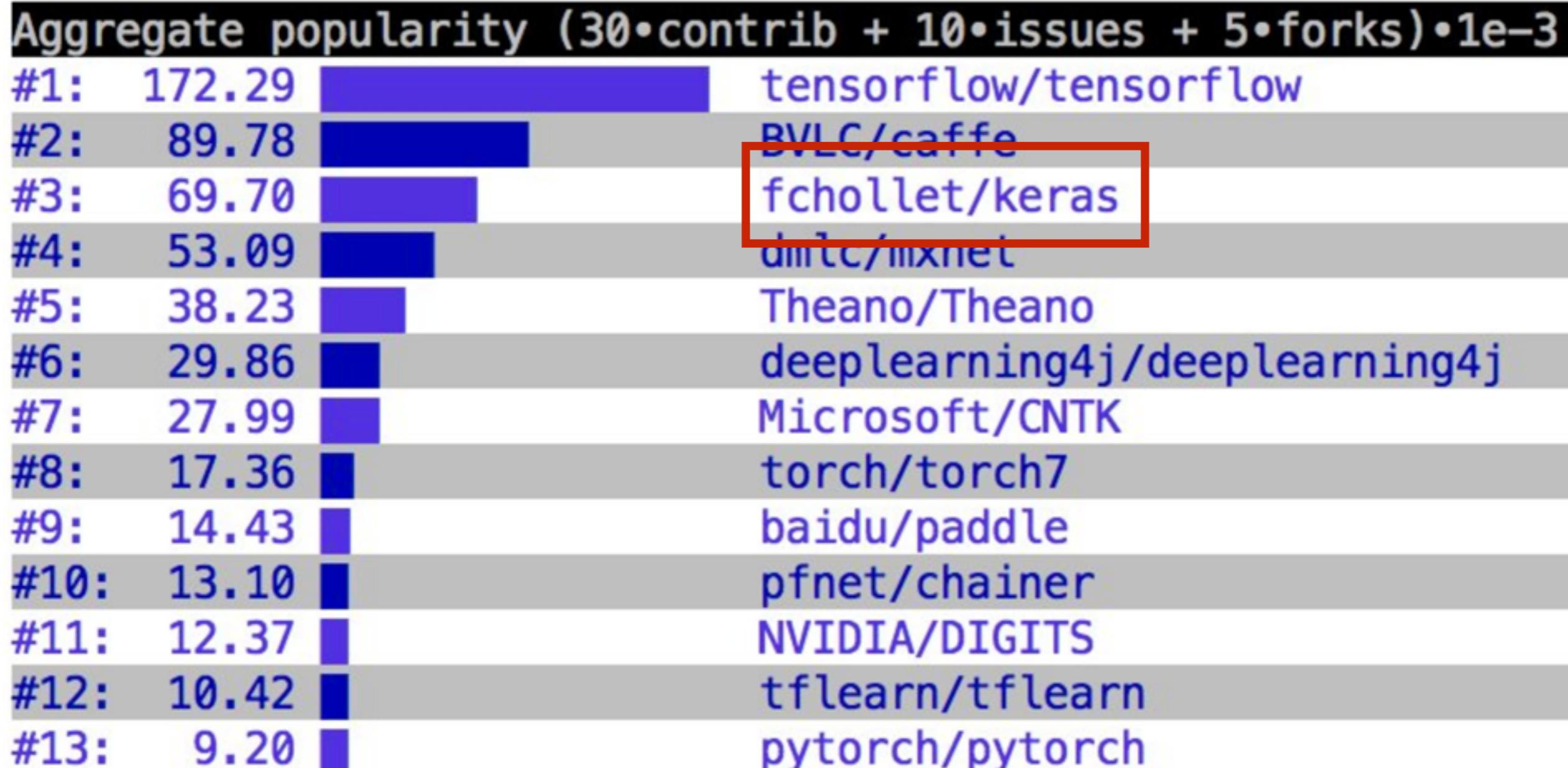
```
correct_pred = tf.equal(tf.argmax(y, 1),  
                      tf.argmax(y_, 1))  
correct_pred_float = tf.cast(correct_pred,  
                             tf.float32)  
acc_fun = tf.reduce_mean(correct_pred_float)  
accuracy = acc_fun.eval(feed_dict={  
    x: mnist.test.images,  
    y_: mnist.test.labels} )  
print accuracy
```

# Evaluate

```
correct_pred = tf.equal(tf.argmax(y, 1),  
                      tf.argmax(y_, 1))  
correct_pred_float = tf.cast(correct_pred,  
                             tf.float32)  
acc_fun = tf.reduce_mean(correct_pred_float)  
accuracy = acc_fun.eval(feed_dict={  
    x: mnist.test.images,  
    y_: mnist.test.labels})  
print accuracy
```



# Deep learning libraries: Accumulated GitHub metrics



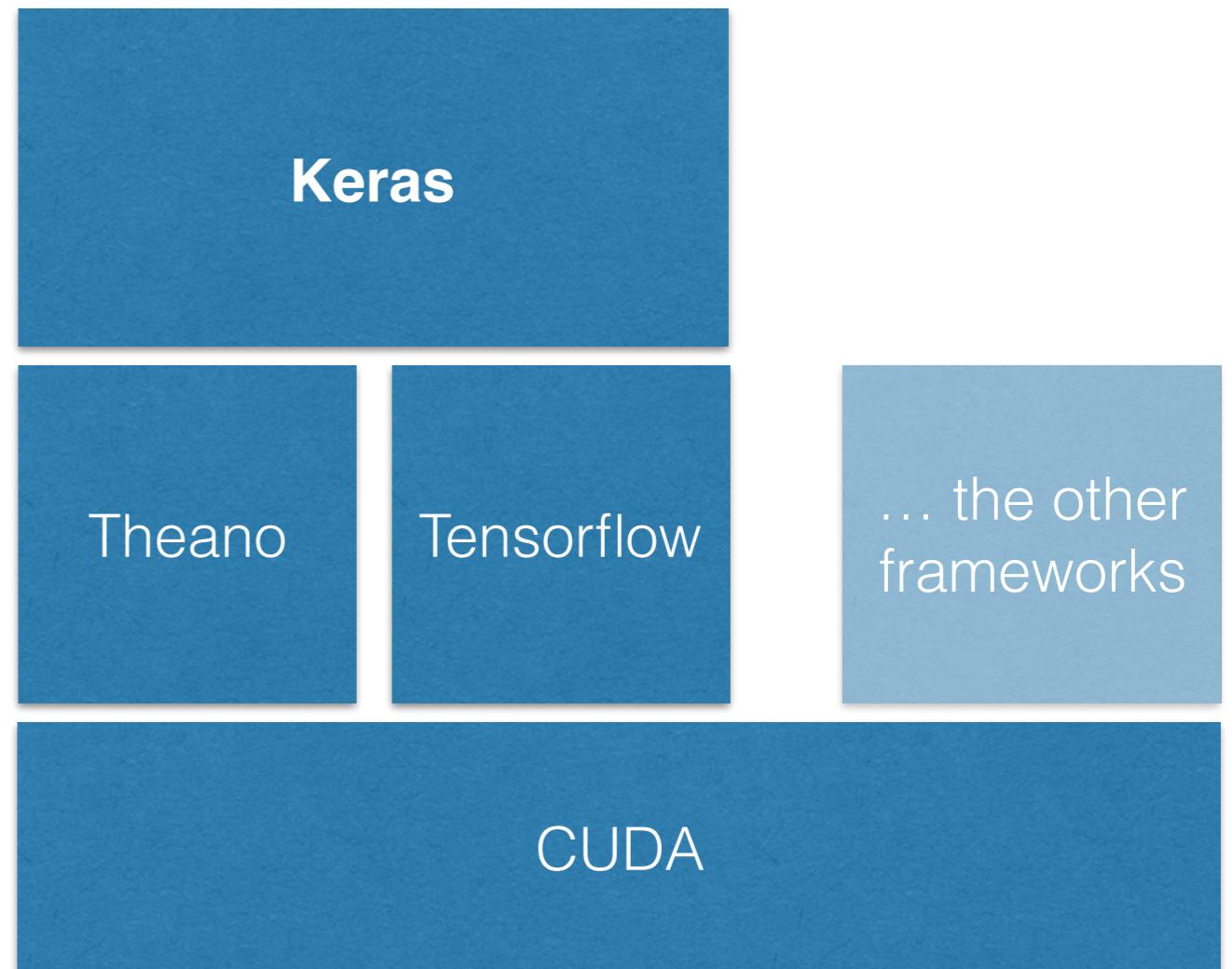
François Chollet @fchollet · Feb 11

Time for an update: what does the deep learning library landscape look like, seen from GitHub?

[pic.twitter.com/QDZyvLrYBd](https://pic.twitter.com/QDZyvLrYBd)

# Keras

- NN framework for fast prototyping
- High level
- It supports Tensorflow AND Theano
- Soon part of the official TensorFlow



# Tensor flow

- Features
  - `speciality` layers
  - most recent loss and activation functions
  - training ‘hacks’ (dropout)
  - helper functions
- Abstraction
- Speed
  - CPU/GPU speed
  - multi GPU support
- Ease of customization
- Base language
  - Python , C++, Java, Rust(!)...

**=> We will use Keras**

# Init

```
TF import tensorflow as tf
sess = tf.InteractiveSession()

K from keras.models import Sequential
from keras.layers.core import Dense, Dropout,
    Activation
from keras.optimizers import SGD
```

# Setup variables

```
TF x = tf.placeholder(tf.float32,  
                     shape=[None, 784])  
y_ = tf.placeholder(tf.float32,  
                     shape=[None, 10])  
W = tf.Variable(tf.zeros([784, 10]))  
b = tf.Variable(tf.zeros([10]))  
  
sess.run(tf.global_variables_initializer())
```

```
K # not necessary
```

# Define the network

```
TF y = tf.matmul(x,W) + b
softmax = tf.nn.\ 
    softmax_cross_entropy_with_logits(y_, y)
cross_entropy = tf.reduce_mean(softmax)
GD = tf.train.GradientDescentOptimizer
train_step = GD().minimize(cross_entropy)
```

```
K model = Sequential()
model.add(Dense(10, input_shape=(784,)))
model.add(Activation('sigmoid'))
model.compile(loss='categorical_crossentropy',
    optimizer=SGD(), metrics=['accuracy'])
```

# Train

```
TF for _ in range(1000):
    batch = mnist.train.next_batch(100)
    train_step.run(feed_dict={x: batch[0],
                                y_: batch[1]})
```

```
K history = model.fit(X_train, Y_train,  
batch_size=batch_size,  
nb_epoch=nb_epoch, verbose=1,  
validation_data=(X_test, Y_test))
```

# Evaluate

TF

```
correct_pred = tf.equal(tf.argmax(y, 1),  
                      tf.argmax(y_, 1))  
correct_pred_float = tf.cast(correct_pred,  
                             tf.float32)  
acc_fun = tf.reduce_mean(correct_pred_float)  
accuracy = acc_fun.eval(feed_dict={  
    x: mnist.test.images,  
    y_: mnist.test.labels})  
print accuracy
```

K

```
score = model.evaluate(X_test, Y_test, verbose=0)  
print 'Test score:', score[0]  
print 'Test accuracy:', score[1]
```

Help, my neural  
network is not learning!

# What ‘Advanced’ techniques do we already talked about?

- Softmax
- Better weight initialization
- L2 Weight Regularization

# Broad strategy

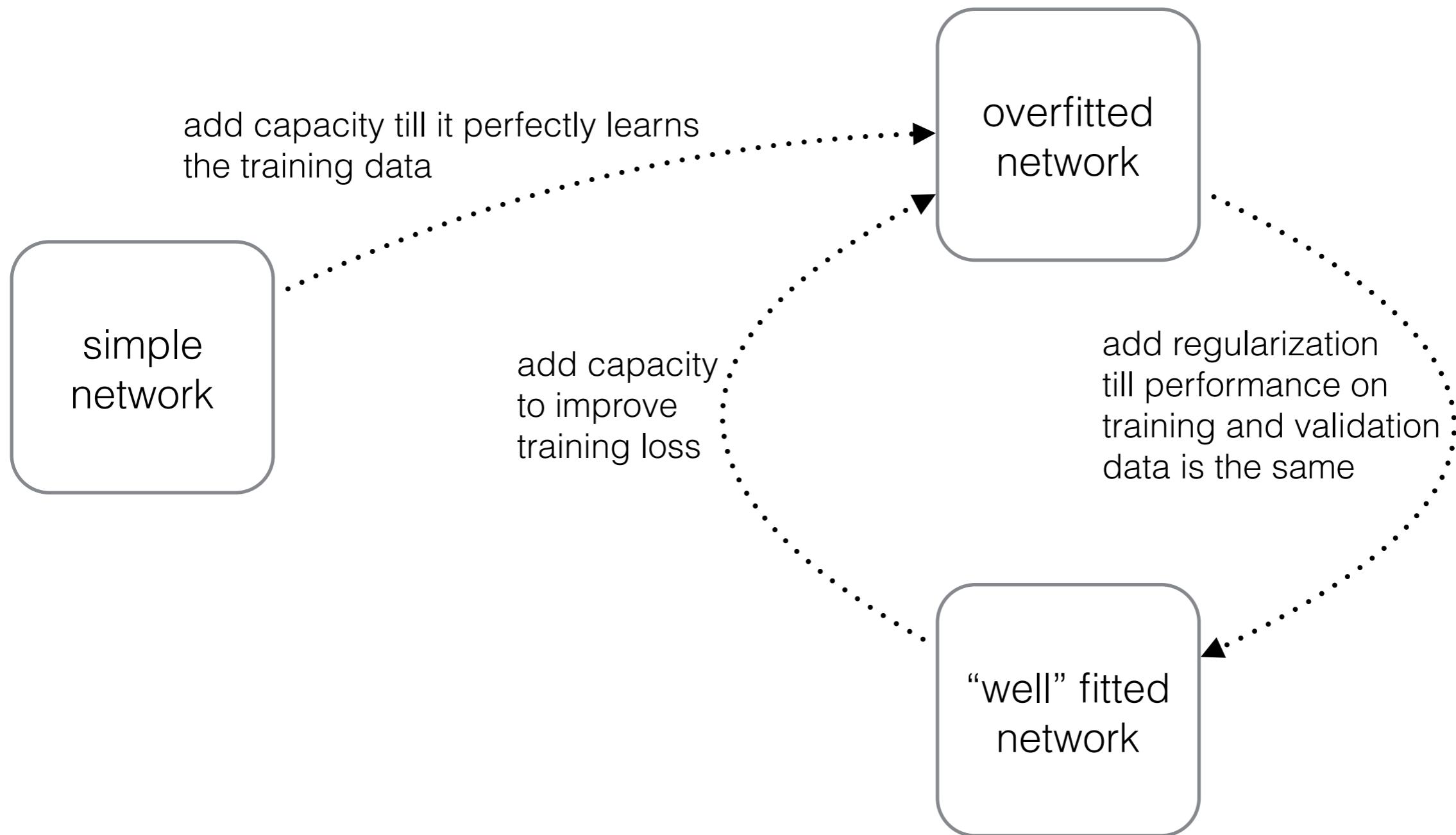
Simplify, simplify, simplify!

- Use smaller network
  - Use subset of classes
  - Monitor progress more often with less data
- Fail fast and try new parameter

# What doesn't make a network simpler?

- outdated activation functions
- ‘wrong’ loss
- vanilla weight initialization

# Regularization and network capacity

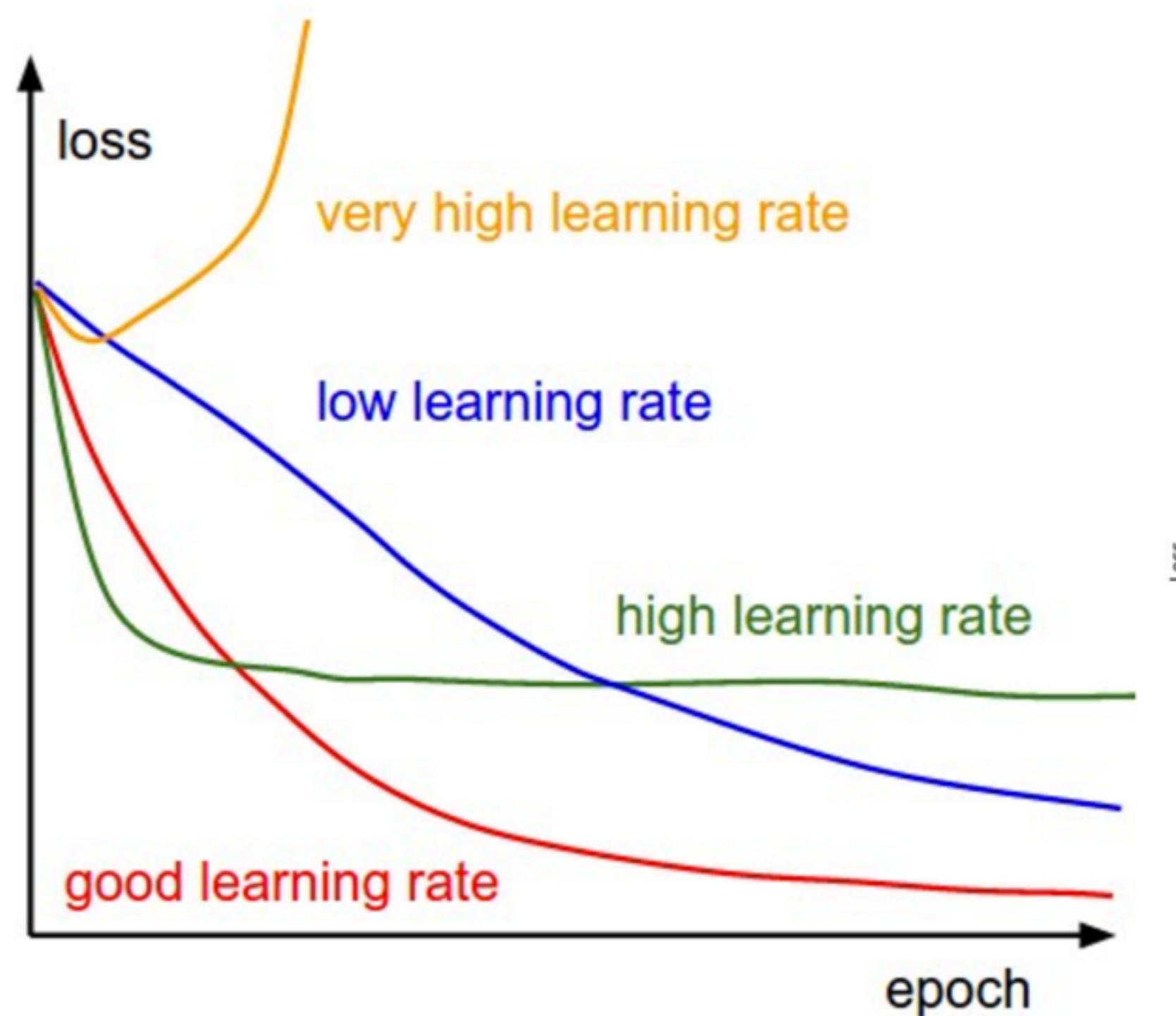


# Learning rate

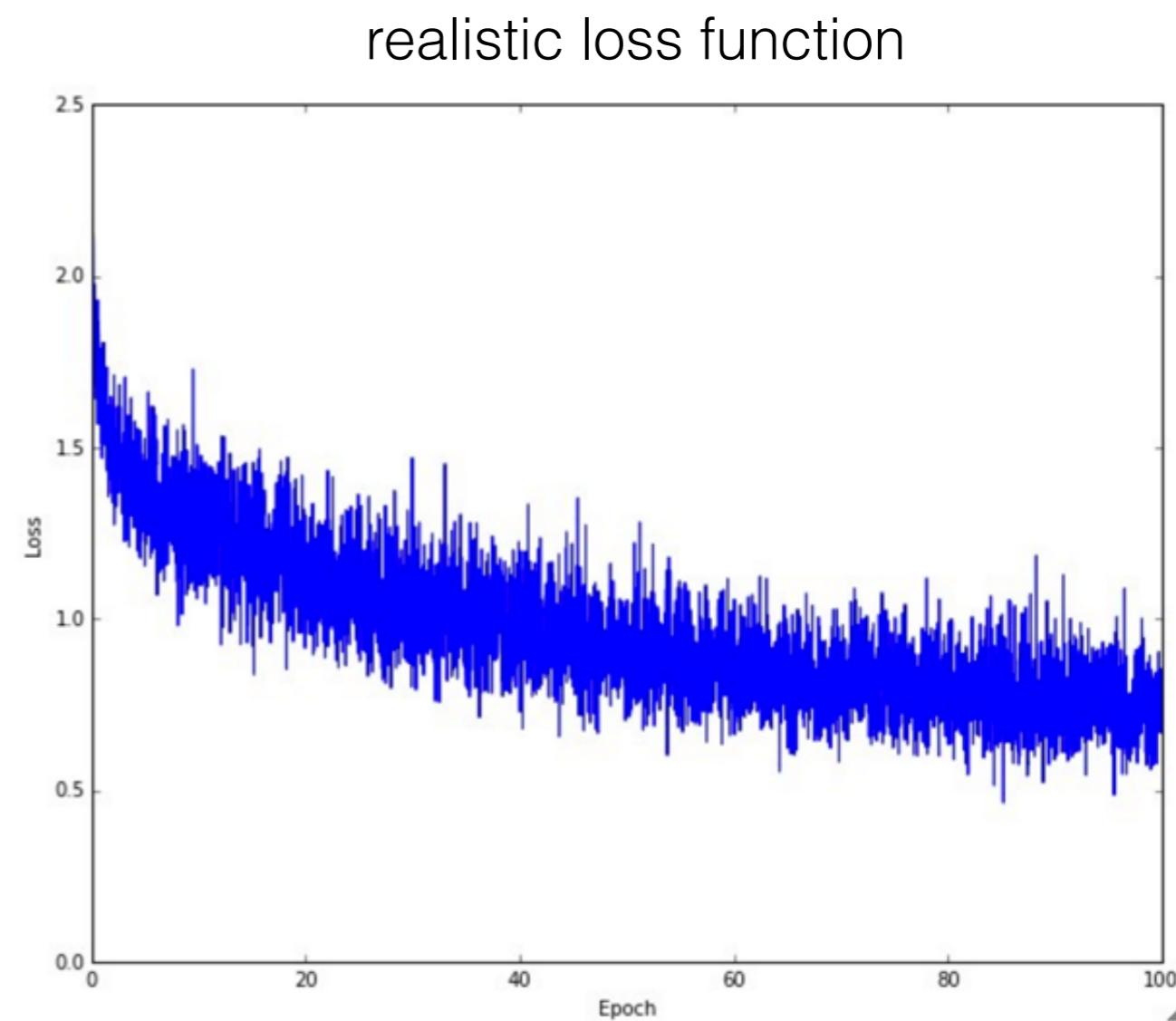
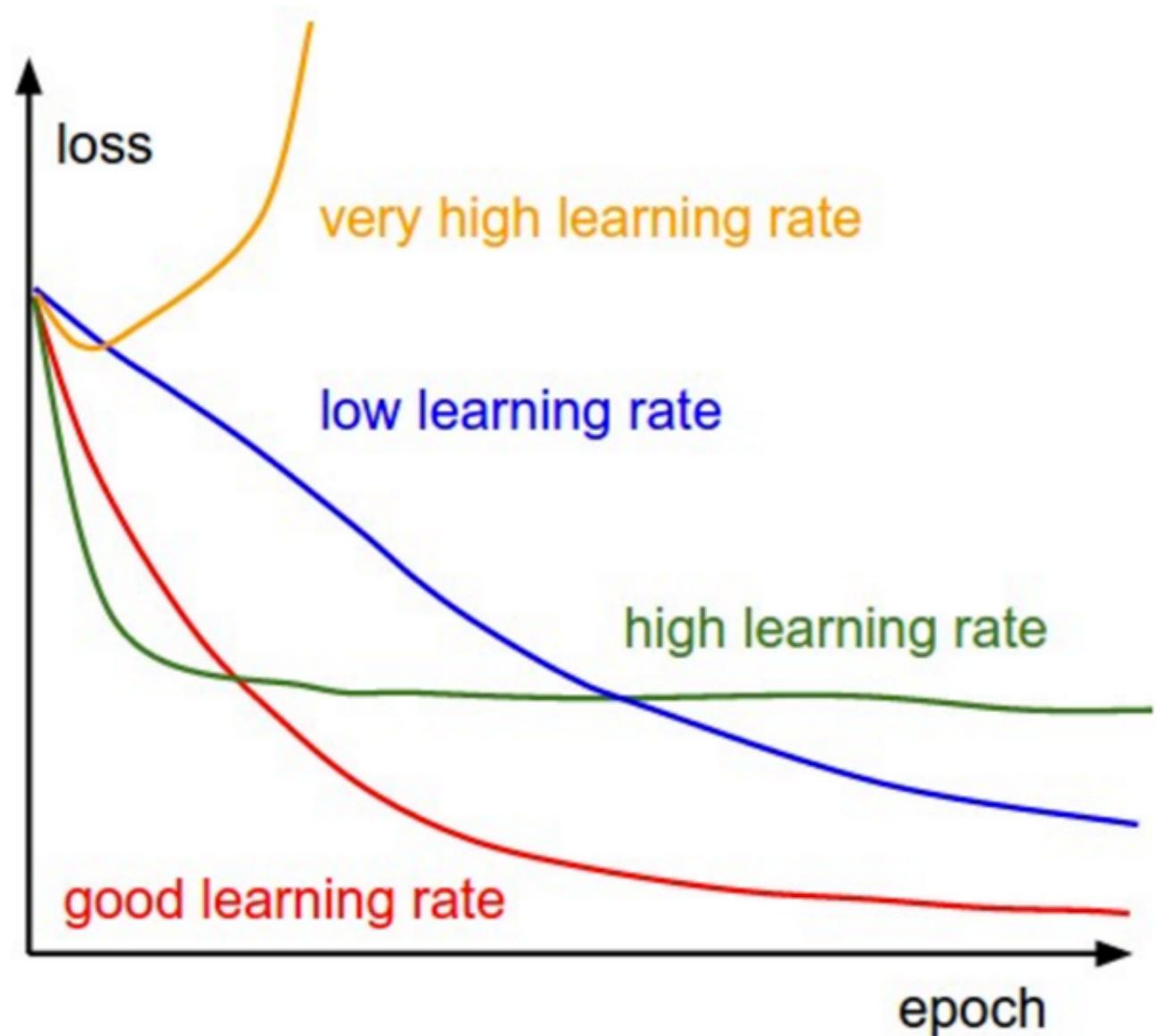
1. Try different magnitudes [1, 0.1, 10, 0.01, 100, ...]
2. Take the biggest learning rate that doesn't 'explode' the network.
3. Half it

Detour: reading the  
loss ‘function’

# Reading the loss ‘function’



# Reading the loss ‘function’



# Mini-batch size

## **Batch size 1**

- also called online learning
- weights are updated often

## **Batch size $1 < M < n_{\text{complete}}$**

- loss is computed on ‘outdated’ weights
- allows to profit of the special optimizations of the hard- and software implementation

**Rule of thumb:** Use the largest batch size, that you can fit into RAM - usually  $>10$  and  $<500$

**To validate:** Plot validation accuracy over time (not epochs) and choose sweet point

# Dependence: Learning rate and batch size

Batch size and learning rate are anti `proportional` in a non-linear way.

Rule of thumb:

**batch\_size \* X => learning\_rate / X (or sqrt(X))**

# Layers

**Fully connected layers:** Keep getting smaller from input to output

**Convolutional layers:** Set filter size 3x3, tune number of filters. Go from few filters to many filters.

**(Max) Pooling layers:** 2x2 or 3x3 and increase number filters

# Data

- needs to be normalized
  - normally [0..1]
- augmentations can be used to increase dataset
  - mirroring, cropping, skewing, rotating...
  - also functions as regularization

# Data

- needs to be normalized
  - normally [0..1]
- augmentations can be used to increase dataset
  - mirroring, cropping, skewing, rotating...
  - also functions as regularization

Why do we need to  
normalize??

# Train MNIST with Keras

Play around with:

- learning rate
- batch size
- network architecture

Try out:

- Softmax
- Better weight initialization
- L2 Weight Regularization
- Augmentation
- See **<https://keras.io/>**

`~/keras_mnist.ipynb`

# Train MNIST with Keras

Play around with:

- learning rate
- batch size
- network architecture

Try out:

- Softmax
- Better weight initialization
- L2 Weight Regularization
- Augmentation
- See <https://keras.io/>

**Set verbose=2  
in model.fit()**

~ / keras\_mnist.ipynb